

Version dated: 2 July, 2015

Accounting for genotype uncertainty in the estimation of allele frequencies in autopolyploids

Paul D. Blischak^{1,*}, Laura S. Kubatko^{1,2} and Andrea D. Wolfe¹

¹*Department of Evolution, Ecology and Organismal Biology, Ohio State University,
318 W. 12th Avenue, Columbus, OH 43210, USA.*

²*Department of Statistics, Ohio State University,
1958 Neil Avenue, Columbus, OH 43210, USA.*

***Corresponding author:** Paul Blischak, Ohio State University, Dept. of Evolution, Ecology and Organismal Biology, 318 W. 12th Avenue, Columbus, OH 43210. E-mail: blischak.4@osu.edu.

Running title: Genotype uncertainty in autopolyploids

Abstract

Despite the ever increasing opportunity to collect large-scale datasets for population genomic analyses, the use of high throughput sequencing to study populations of polyploids has seen little application. This is due in large part to problems associated with determining allele copy number in the genotypes of polyploid individuals (allelic dosage uncertainty–ADU), which complicates the calculation of important quantities such as allele frequencies. This well-known problem has hindered population genetic studies in polyploids even though various solutions to circumvent the difficulty of estimating polyploid genotypes have been proposed. Additional complications arise because of the mixed inheritance patterns and variable reproductive modes that are characteristic of many polyploid taxa, making the development of population genetic models for polyploids especially difficult. Here we describe a statistical model to estimate biallelic SNP frequencies in a population of autopolyploids using high throughput sequencing data in the form of read counts. Uncertainty in the number of copies of an allele in an individual’s genotype is accounted for by treating genotypes as a latent variable in a hierarchical Bayesian model. In this way, we bridge the gap from data collection (using techniques such as restriction-site associated DNA sequencing) to allele frequency estimation in a unified inferential framework by summing over genotype uncertainty. Simulated datasets were

28 generated under various conditions for both tetraploid and hexaploid populations to evaluate the
29 model's performance and to help guide the collection of empirical data. We also discuss potential
30 extensions to generalize our model and its application to polyploids.

31 (**Keywords:** allelic dosage uncertainty, allele frequencies, hierarchical Bayesian modeling, polyploidy,
32 population genomics)

33 Introduction

34 Biologists have long been fascinated by the occurrence of whole genome duplication (WGD) in
35 natural populations and have recognized its role in the generation of biodiversity (Clausen *et al.* 1940;
36 Stebbins 1950; Grant 1971; Otto & Whitton 2000). Though WGD is thought to have occurred at some
37 point in nearly every branch of the Tree of Life (plants, animals and fungi), it is a particularly common
38 phenomenon in plants and is regarded by many to be an important factor in plant diversification
39 (Wood *et al.* 2009; Soltis *et al.* 2009; Scarpino *et al.* 2014). The role of polyploidy in plant evolution
40 was originally considered by some to be a “dead-end” (Stebbins 1950; Wagner 1970; Soltis *et al.*
41 2014) but, since its first discovery in the early twentieth century, polyploidy has been continually
42 studied in nearly all areas of botany (Winge 1917; Winkler 1916; Clausen *et al.* 1945; Grant 1971;
43 Stebbins 1950; Soltis *et al.* 2003, 2010; Soltis & Soltis 2009; Ramsey & Ramsey 2014). Though fewer
44 examples of WGD are currently known for animal systems, groups such as amphibians, fish, and
45 reptiles all exhibit polyploidy (Allendorf & Thorgaard 1984; Gregory & Mable 2005). Ancient genome
46 duplications are also thought to have played an important role in the evolution of both plants and
47 animals, occurring in the lineages preceding the seed plants, angiosperms and vertebrates (Ohno
48 1970; Otto & Whitton 2000; Furlong & Holland 2001; Jiao *et al.* 2011). These ancient WGD events
49 during the early history of seed plants and angiosperms have been followed by several more WGDs in
50 all major plant groups (Cui *et al.* 2006; Scarpino *et al.* 2014; Cannon *et al.* 2014). Recent experimental
51 evidence has also demonstrated increased survivorship and adaptability to foreign environments of
52 polyploid taxa when compared with their lower ploidy relatives (Ramsey 2011; Selmecki *et al.* 2015).

53 Polyploids are generally divided into two types based on how they are formed: auto- and
54 allopolyploids. Autopolyploids form when a WGD event occurs within a single evolutionary lineage
55 and typically have polysomic inheritance. Allopolyploids are formed by hybridization between two
56 separately evolving lineages followed by WGD and are thought to have mostly disomic inheritance.
57 Multivalent chromosome pairing during meiosis can occur in allopolyploids, however, resulting in
58 mixed inheritance patterns across loci in the genome [segmental allopolyploids] (Stebbins 1950).
59 Autopolyploids can also undergo double reduction, a product of multivalent chromosome pairing
60 wherein segments from sister chromatids move together during meiosis—resulting in allelic inheritance
61 that breaks away from a strict pattern of polysomy (Haldane 1930). Autopolyploidy was also thought
62 to be far less common than allopolyploidy, but recent studies have concluded that autopolyploidy
63 occurs much more frequently than originally proposed (Soltis *et al.* 2007; Parisod *et al.* 2010).

64 The theoretical treatment of population genetic models in polyploids has its origins in the Modern
65 Synthesis with Fisher, Haldane and Wright each contributing to the development of some of the
66 earliest mathematical models for understanding the genetic patterns of inheritance in polyploids.
67 Among the first of these works was Haldane's 1930 paper on autopolyploid inheritance in $2k$ -ploid

68 ($k = 2, 3, \dots$) organisms. Influenced in part by the works of Hermann J. Muller in tetraploid species of
69 *Primula* (1914) and W. J. C. Lawrence in octoploid species of *Dahlia* (1929), Haldane generalized the
70 combinatorial formulas for determining the frequencies of the different possible gametes formed from
71 all genotype combinations for a $2k$ -ploid. He also considered additional factors influencing gamete
72 frequencies such as double reduction and the effects of partial selfing (Haldane 1930). Fisher's interest
73 in polyploidy stemmed largely from observations made in the plant genus *Lythrum*, which exhibited
74 conspicuous patterns of trimorphic heterostyly (Fisher 1941). Empirical works by Nora Barlow (1913,
75 1923), as well as initial investigations into the inheritance patterns of the three style types (Short, Mid,
76 Long) by E. M. East (1927) formed the basis for Fisher's formulation of a model for the inheritance
77 patterns of the Mid length style form in *Lythrum salicaria* (Fisher 1941). He later added to this work
78 by considering double reduction in the inheritance of the Mid length style and complemented his
79 theoretical work through a collaboration with Kenneth Mather to complete crossing experiments
80 (Fisher 1943; Fisher & Mather 1943). Wright's contributions were concerned with the calculation of
81 the distribution of allele frequencies in a $2k$ -ploid and were largely an extension of his classic 1931
82 paper, *Evolution in Mendelian populations*, and a previously published manuscript describing similar
83 processes in diploids (Wright 1931, 1937, 1938). Wright was among the first to consider mutation,
84 migration, selection and inbreeding in his formulation of the distribution of gene frequencies, which
85 helped to establish future ideas about modeling allelic diffusion in a population. For example, it was
86 noted by Kimura (1964) that much of the work on diffusion equations in population genetics could be
87 applied to polyploids in a manner similar to Wright's derivation of the allele frequency distribution
88 in polyploids.

89 The foundation laid down by these early papers has led to the continuing development of
90 population genetic models for polyploids, including models for understanding the rate of loss of
91 genetic diversity and extensions of the coalescent in autotetraploids, as well as modifications of the
92 multispecies coalescent for the inference of species networks containing allotetraploids (Moody *et al.*
93 1993; Arnold *et al.* 2012; Jones *et al.* 2013). Much of this progress was described in a review by
94 Dufresne *et al.* (2014), who outlined the current state of population genetics in polyploids regarding
95 both molecular techniques and statistical models. Not surprisingly, one of the most promising
96 developments for the future of population genetics in polyploids is the advancement of sequencing
97 technologies. A particularly common method of gathering large datasets for genome scale inferences
98 is restriction-site associated DNA sequencing [RADseq] (Miller *et al.* 2007; Baird *et al.* 2008; Puritz
99 *et al.* 2014). However, despite its popularity for population genetic inferences at the diploid level,
100 there are many fewer examples of RADseq experiments conducted on polyploid taxa (but see Ogden
101 *et al.* 2013; Wang *et al.* 2013; Logan-Young *et al.* 2015). Among the primary reasons for the dearth
102 in applying RADseq to polyploids is the issue of allelic dosage uncertainty (ADU), or the inability to
103 fully determine the genotype of a polyploid organism when it is partially heterozygous at a given
104 locus. This is the same problem that has been encountered by other codominant markers such as
105 microsatellites, which have been commonly used for population genetic analyses in polyploids. One
106 way of dealing with allelic dosage that has been used for multi-allelic microsatellite markers has been
107 to code alleles as either present or absent based on electropherogram readings (allelic phenotypes)
108 and to analyze the resulting dominant data using a program such as POLYSAT (Clark & Jasieniuk
109 2011; Dufresne *et al.* 2014). de Silva *et al.* (2005) developed a method for inferring allele frequencies
110 using observed allelic phenotype data and used an expectation-maximization algorithm to deal
111 with the incomplete genotype data resulting from ADU. Attempts to directly infer the genotype
112 of polyploid microsatellite loci have also been successfully completed in some cases by using the

113 relative electropherogram peak heights of the alleles in the genotypes (Esselink *et al.* 2004). The
114 estimation problem would be similar for biallelic SNP data collected using RADseq, where a partially
115 heterozygous polyploid will have high throughput sequencing reads containing both alleles. For a
116 tetraploid, the possible genotypes for a partial heterozygote (alleles A and B) would be AAAB, AABB
117 and ABBB. For a hexaploid they are AAAAAB, AAAABB, AAABBB, AABBBB and ABBBBB. In
118 general, the number of possible genotypes for a biallelic locus of a partially heterozygous K-ploid
119 ($K = 3, 4, 5, \dots$) is $K - 1$. A possible solution to this problem for SNPs would be to try to use
120 existing genotype callers and to rely on the relative number of sequencing reads containing the two
121 alleles (similar to what was done for microsatellites). However, this could lead to erroneous inferences
122 when genotypes are simply fixed at point estimates based on read proportions without considering
123 estimation error. Furthermore, when sequencing coverage is low, the number of genotypes that will
124 appear to be equally probable increases with ploidy, making it difficult to distinguish among the
125 possible partially heterozygous genotypes.

126 In this paper we describe a model that aims to address the problems associated with ADU by
127 treating genotypes as a latent variable in a hierarchical Bayesian model and using high throughput
128 sequencing read counts as data (Figure 1). In this way we preserve the uncertainty that is inherent
129 in the genotypes of partially heterozygous polyploids by inferring a probability distribution across all
130 possible values of the genotype, rather than treating the genotypes as being directly observed. This
131 approach has been used by Buerkle & Gompert (2013) to deal with uncertainty in calling genotypes
132 in diploids and the work we present here builds off of their earlier models. Our model assumes that
133 the ploidy level of the population is known and that the genotypes of individuals in the population
134 are drawn from a single underlying allele frequency for each locus. These assumptions imply that
135 alleles in the population are undergoing polysomic inheritance without double reduction, which most
136 closely adheres to the inheritance patterns of an autopolyploid. We acknowledge that the model in its
137 current form is an oversimplification of biological reality and realize that it does not apply to a large
138 portion of polyploid taxa. Nevertheless, we believe that accounting for ADU by modeling genotype
139 uncertainty has the potential to be applied more broadly via modifications of the probability model
140 used for the inheritance of alleles, which could lead to more generalized population genetic models
141 for polyploids (see the **Extensibility** section of the **Discussion**).

142 Materials and Methods

143 Our goal is to estimate the frequency of a reference allele for each locus sampled from a population
144 of known ploidy (ψ), where the reference allele can be chosen arbitrarily between the two alleles at
145 a given biallelic SNP. To do this we extend the population genomic models of Buerkle & Gompert
146 (2013), which employ a Bayesian framework to model high throughput sequencing reads (\mathbf{T} , \mathbf{R}),
147 genotypes (\mathbf{G}) and allele frequencies (\mathbf{p}), to the case of arbitrary ploidy. The idea behind the model
148 is to view the sequencing reads gathered for an individual as a random sample from the unobserved
149 genotype at each locus. Genotypes can then be treated as a parameter in a probability model that
150 governs how likely it is that we see a particular number of sequencing reads carrying the reference
151 allele. Similarly, we can treat genotypes as a random sample from the underlying allele frequency
152 in the population (assuming Hardy-Weinberg equilibrium). For our model, a genotype is simply a
153 count of the number of reference alleles at a locus which can range from 0 (a homozygote with no
154 reference alleles in the genotype) to ψ (a homozygote with only reference alleles in the genotype).
155 All whole numbers in between 0 and ψ represent partially heterozygous genotypes. This hierarchical

156 setup addresses the problems associated with ADU by treating genotypes as a latent variable that
 157 can be integrated out using Markov chain Monte Carlo (MCMC).

158 Model setup

159 Here we consider a sample of N individuals from a single population of ploidy level ψ sequenced at L
 160 unlinked SNPs. The data for the model consist of two matrices containing counts of high throughput
 161 sequencing reads mapping to each locus for each individual: \mathbf{R} and \mathbf{T} . The $N \times L$ matrix \mathbf{T} contains
 162 the total number of reads sampled at each locus for each individual. Similarly, \mathbf{R} is an $N \times L$ matrix
 163 containing the number of sampled reads with the reference allele at each locus for each individual.
 164 Then for individual i at locus ℓ , we model the number of sequencing reads containing the reference
 165 allele ($r_{i\ell}$) as a Binomial random variable conditional on the total number of sequencing reads ($t_{i\ell}$),
 166 the underlying genotype ($g_{i\ell}$) and a constant level of sequencing error (ϵ)

$$P(r_{i\ell}|t_{i\ell}, g_{i\ell}, \epsilon) = \binom{t_{i\ell}}{r_{i\ell}} \begin{cases} \epsilon^{r_{i\ell}}(1-\epsilon)^{t_{i\ell}-r_{i\ell}} & \text{if } g_{i\ell} = 0, \\ \left(\frac{g_{i\ell}}{\psi}\right)^{r_{i\ell}} \left(1 - \frac{g_{i\ell}}{\psi}\right)^{t_{i\ell}-r_{i\ell}} & \text{if } g_{i\ell} = 1, \dots, \psi - 1, \\ (1-\epsilon)^{r_{i\ell}} \epsilon^{t_{i\ell}-r_{i\ell}} & \text{if } g_{i\ell} = \psi. \end{cases} \quad (1)$$

167 If we assume conditional independence of the sequencing reads given the genotypes, the joint
 168 probability distribution for sequencing reads is given by

$$P(\mathbf{R}|\mathbf{T}, \mathbf{G}, \epsilon) = \prod_{\ell=1}^L \prod_{i=1}^N P(r_{i\ell}|t_{i\ell}, g_{i\ell}, \epsilon). \quad (2)$$

169 Since the $r_{i\ell}$'s are the data that we observe, the product of $P(r_{i\ell}|t_{i\ell}, g_{i\ell}, \epsilon)$ across loci and individuals
 170 will form the likelihood in the model. An important consideration here is that when $g_{i\ell}$ is equal to
 171 0 or ψ (i.e., when the genotype is homozygous) but the sequence data collected for the individual
 172 contain both alleles, the likelihood will be 0. To correct for this, we include error (ϵ) into the model.
 173 The intuition behind including error is that reads sampled from a locus that is truly homozygous
 174 can still contain more than one allele due to sequencing errors, giving the false impression that the
 175 individual may be a partial heterozygote. If we assume that the probability of a sequencing error is
 176 smaller than the probability of being truly heterozygous (i.e., containing at least one reference allele
 177 in the genotype), then the probability model we use for reference reads should be able to distinguish
 178 between homozygotes with reads containing sequencing errors and partial heterozygotes. When we
 179 do this, the probability distribution for $r_{i\ell}$ given $g_{i\ell}$ and ϵ is split into $\psi + 1$ cases as above in Eq. 1.

180 The next level in the hierarchy is the conditional prior for genotypes. We model each $g_{i\ell}$ as a
 181 Binomial random variable conditional on the ploidy level of the population and the frequency of the
 182 reference allele for locus ℓ (p_ℓ):

$$P(g_{i\ell}|\psi, p_\ell) = \binom{\psi}{g_{i\ell}} p_\ell^{g_{i\ell}} (1-p_\ell)^{\psi-g_{i\ell}}.$$

183 We also assume that the genotypes of the sampled individuals are conditionally independent given the
 184 allele frequencies, which is equivalent to taking a random sample from a population in Hardy-Weinberg
 185 equilibrium. Factoring the distribution for genotypes and taking the product across loci and individuals

186 gives us the joint probability distribution of genotypes given the ploidy level of the population and
 187 the vector of allele frequencies at each locus ($\mathbf{p} = \{p_1, \dots, p_L\}$):

$$P(\mathbf{G}|\psi, \mathbf{p}) = \prod_{\ell=1}^L \prod_{i=1}^N P(g_{i\ell}|\psi, p_{\ell}). \quad (3)$$

188 The final level of the model is the prior distribution on allele frequencies. Assuming *a priori*
 189 independence across loci, we use a Beta distribution with parameters α and β both equal to 1 as our
 190 prior distribution for each locus. A Beta(1,1) is equivalent to a Uniform distribution over the interval
 191 $[0, 1]$, making our choice of prior uninformative. The joint posterior distribution of allele frequencies
 192 and genotypes is then equal to the product across all loci and all individuals of the likelihood, the
 193 conditional prior on genotypes and the prior distribution on allele frequencies up to a constant of
 194 proportionality

$$\begin{aligned} P(\mathbf{p}, \mathbf{G}|\mathbf{R}, \epsilon) &\propto P(\mathbf{R}|\mathbf{T}, \mathbf{G}, \epsilon)P(\mathbf{G}|\psi, \mathbf{p})P(\mathbf{p}) \\ &= \prod_{\ell=1}^L \prod_{i=1}^N P(r_{i\ell}|t_{i\ell}, g_{i\ell}, \epsilon)P(g_{i\ell}|\psi, p_{\ell})P(p_{\ell}). \end{aligned} \quad (4)$$

195 The marginal posterior distribution for allele frequencies can be obtained by summing over genotypes

$$P(\mathbf{p}|\mathbf{R}, \epsilon) \propto \sum_{\mathbf{G}} P(\mathbf{p}, \mathbf{G}|\mathbf{R}, \epsilon). \quad (5)$$

196 It would also be possible to examine the marginal posterior distribution of genotypes but here we will
 197 only focus on allele frequencies (see **Discussion** for potential applications of the marginal distribution
 198 of genotypes).

199 Full conditionals and MCMC using Gibbs sampling

200 We estimate the joint posterior distribution for allele frequencies and genotypes in Eq. 4 using MCMC.
 201 This is done using Gibbs sampling of the states (\mathbf{p}, \mathbf{G}) in a Markov chain by alternating samples
 202 from the full conditional distributions of \mathbf{p} and \mathbf{G} . Given the setup for our model using Binomial
 203 and Beta distributions (which form a conjugate family), analytical solutions for these distributions
 204 can be readily acquired (Gelman *et al.* 2014). The full conditional distribution for allele frequencies
 205 is Beta distributed and is given by Eq. 6 below:

$$p_{\ell} | g_{i\ell}, r_{i\ell}, \epsilon \sim \text{Beta} \left(\alpha = \sum_{i=1}^N g_{i\ell} + 1, \beta = \sum_{i=1}^N (\psi - g_{i\ell}) + 1 \right), \quad \text{for } \ell = 1, \dots, L. \quad (6)$$

206 This full conditional distribution for p_{ℓ} has a natural interpretation as it is roughly centered at the
 207 proportion of sampled alleles carrying the reference allele divided by the total number of alleles
 208 sampled given the current state of \mathbf{G} in the Markov chain. The “+1” comes from the prior distribution
 209 and will not have a strong influence on the posterior when the sample size is large.

210 The full conditional distribution for genotypes is split into $\psi + 1$ cases (similar to the conditional
 211 prior), making it a discrete categorical distribution over the possible values for the genotypes $(0, \dots, \psi)$.
 212 Using k as a generic index, the distribution for individual i at locus ℓ is

$$P(g_{i\ell}|g_{(-i)\ell}, p_\ell, r_{i\ell}, \epsilon) = \frac{1}{\mathcal{C}_{i\ell}} \begin{cases} \epsilon^{r_{i\ell}}(1-\epsilon)^{t_{i\ell}-r_{i\ell}}(1-p_\ell)^\psi & \text{for } k = 0, \\ \left(\frac{k}{\psi}\right)^{r_{i\ell}} \left(1-\frac{k}{\psi}\right)^{t_{i\ell}-r_{i\ell}} \binom{\psi}{k} p_\ell^k (1-p_\ell)^{\psi-k} & \text{for } k = 1, \dots, \psi - 1, \\ (1-\epsilon)^{r_{i\ell}} \epsilon^{t_{i\ell}-r_{i\ell}} p_\ell^\psi & \text{for } k = \psi, \end{cases} \quad (7)$$

213 where $g_{(-i)\ell}$ is the value of the genotypes for all sampled individuals excluding individual i and $\mathcal{C}_{i\ell}$ is
 214 a normalizing constant equal to the sum of all of the terms:

$$\mathcal{C}_{i\ell} = \epsilon^{r_{i\ell}}(1-\epsilon)^{t_{i\ell}-r_{i\ell}}(1-p_\ell)^\psi + (1-\epsilon)^{r_{i\ell}} \epsilon^{t_{i\ell}-r_{i\ell}} p_\ell^\psi + \sum_{k=1}^{\psi-1} \left(\left(\frac{k}{\psi}\right)^{r_{i\ell}} \left(1-\frac{k}{\psi}\right)^{t_{i\ell}-r_{i\ell}} \binom{\psi}{k} p_\ell^k (1-p_\ell)^{\psi-k} \right).$$

215 The full conditional distribution for genotypes can be seen as the product of two quantities: (1)
 216 the probability of each of the possible genotypes based on the observed reference reads and (2)
 217 the probability of drawing each genotype value based on the current value for the frequency of the
 218 reference allele for locus ℓ in the population.

219 We begin our Gibbs sampling algorithm in a random position in parameter space through the use
 220 of uniform probability distributions. The genotype matrix is initialized with random draws from a
 221 Discrete Uniform distribution from 0 to ψ and the initial allele frequencies are drawn from a Uniform
 222 distribution on the interval $[0,1]$.

223 Simulation study

224 Simulations were performed to assess error rates in allele frequency estimation for tetraploid and
 225 hexaploid populations ($\psi = 4$ and 6, respectively). Data were generated under the model by
 226 sampling genotypes from a Binomial distribution conditional on a fixed, known allele frequency
 227 ($p_\ell = 0.01, 0.05, 0.1, 0.2, 0.4$). Total read counts were simulated for a single locus using a Poisson
 228 distribution with mean coverage equal to 5, 10, 20, 50 or 100 reads per individual. We then sampled
 229 the number of sequencing reads containing the reference allele from a Binomial distribution conditional
 230 on the number of total reads, the genotype and sequencing error (Eq. 1; ϵ fixed to 0.01). Finally,
 231 we varied the number of individuals sampled per population ($N = 5, 10, 20, 30$) and ran all possible
 232 combinations of the simulation settings. Each combination of sequencing coverage, individuals sampled
 233 and allele frequency was analyzed using 100 replicates for both tetraploid and hexaploid populations
 234 for a total of 20,000 simulation runs. MCMC analyses using Gibbs sampling were run for 100,000
 235 generations with parameter values stored every 100 samples. The first 25% of the posterior was
 236 discarded as burn-in, resulting in 750 posterior samples for each replicate. Convergence on the
 237 stationary distribution, $P(\mathbf{p}, \mathbf{G}|\mathbf{R}, \epsilon)$, was assessed by examining trace plots for a subset of runs for
 238 each combination of settings and ensuring that the effective sample sizes (ESS) were greater than 200.
 239 Deviations from the known underlying allele frequency used to simulate each data set were assessed
 240 by calculating the posterior mean of each replicate, followed by subtracting the known value from
 241 the calculated means and then calculating their overall standard deviation.

242 All simulations were performed using the R statistical package (R Core Team 2014) on the
243 Oakley cluster at the Ohio Supercomputer Center (<https://osc.edu>). Figures were generated using
244 the R packages GGPlot2 (Wickham 2009), RESHAPE (Wickham 2011) and PLYR (Wickham 2007)
245 with additional figure manipulation completed using Inkscape (<https://inkscape.org>). MCMC
246 diagnostics were done using the CODA package (Plummer *et al.* 2006). All code is available on GitHub
247 (<https://github.com/pblischak/polyfreqs-ms-data>).

248 Results

249 Gibbs sampler

250 Our Gibbs sampling algorithm was able to accurately estimate allele frequencies for a number of
251 simulation settings while simultaneously allowing for genotype uncertainty. There were no indications
252 of a lack of convergence (ESS values > 200) for any of the simulation replicates and all trace plots
253 examined also indicated that the Markov chain had reached stationarity. We have aggregated the
254 scripts for our Gibbs sampler as a developmental R package—POLYFREQS—which is available on
255 GitHub (<https://github.com/pblischak/polyfreqs>). Though POLYFREQS is written in R, it deals
256 with the large datasets that are generated by high throughput sequencing platforms in two ways.
257 First, it takes advantage of R's ability to incorporate C++ code via the RCPP and RCPPARMADILLO
258 packages, allowing for a faster implementation of our MCMC algorithm (Eddelbuettel & François 2011;
259 Eddelbuettel 2013; Eddelbuettel & Sanderson 2014). Second, since the model assumes independence
260 among loci, POLYFREQS can facilitate the process of parallelizing analyses by splitting the total read
261 count and reference read count matrices into subsets of loci which can be analyzed at the same time
262 on separate nodes of a computing cluster.

263 Simulation study

264 Increasing the number of individuals sampled had the largest effect on the accuracy of allele frequency
265 estimation (Figures 2 & 3). Since allele frequencies are population level parameters, it is not surprising
266 that sampling more individuals from the population leads to better estimates. This appears to be the
267 case even when sequencing coverage is quite low (5x, 10x), which corroborates the observations made
268 by Buerkle & Gompert (2013). Lower sequence coverage does affect the posterior distribution for
269 allele frequencies even when the number of individuals sequenced is large however, by increasing the
270 posterior standard deviation (Figure 4). An interesting pattern that emerged during the simulation
271 study is the observation that the allele frequencies closer to 0.5 tend to have higher error rates,
272 which is to be expected given that the variance of a Binomial random variable is highest when the
273 probability of success is 0.5.

274 Discussion

275 The inference of population genetic parameters and the demographic history of non-model polyploid
276 organisms has consistently lagged behind that of diploids. The difficulties associated with these
277 inferences present themselves at two levels. The first of these is the widely known inability to determine

278 the genotypes of polyploids due to ADU. Even though there have been theoretical developments in
279 the description of models for polyploid taxa as early as the 1930s, a large portion of this population
280 genetic theory relies on knowledge about individuals' genotypes (e.g., Haldane 1930; Wright 1938).
281 The second complicating factor is the complexity of inheritance patterns and changes in mating
282 systems that often accompany WGD events. Polyploid organisms can sometimes mate by both
283 outcrossing or selfing, and can display mixed inheritance patterns at different loci in the genome
284 (Dufresne *et al.* 2014). If genotypes were known, then it might be easier to develop and test models
285 for dealing with and inferring rates of selfing versus outcrossing, as well as understanding inheritance
286 patterns across the genome. However, ADU only compounds the problems associated with these
287 inferences, making the development and application of appropriate models far more difficult (but
288 see list of software in Dufresne *et al.* 2014). The model we have presented here deals with the first
289 of these two issues by not treating genotypes as observed quantities. Almost all other methods of
290 genotype estimation for polyploids treat the genotype as the primary parameter of interest. Our
291 model is different in that we still use the read counts generated by high throughput sequencing
292 platforms as our observed data but instead integrate across genotype uncertainty when inferring
293 other parameters, thus bypassing the problems caused by ADU.

294 Despite our focus on bypassing ADU, an important consideration for the model we present here
295 is that, because it approximates the joint posterior distribution of allele frequencies and genotypes, it
296 would also be possible to use the marginal posterior distribution of genotypes to make inferences
297 using existing methods. This could be done using the posterior mode as a maximum *a posteriori*
298 (MAP) estimate of the genotype for downstream analyses, followed by analyzing random samples
299 from the posterior distribution of genotypes. The resulting set of estimates would not constitute a
300 "true" posterior distribution of downstream parameters but would allow researchers to interpret their
301 results based on the MAP estimate of the genotypes while still getting a sense for the amount of
302 variation in their estimates. Using the posterior distribution of genotypes in this way could technically
303 be applied to any type of polyploid, but is only really appropriate for autopolyploids due to the
304 model of inheritance that is used. Other methods for estimating SNP genotypes from high throughput
305 sequencing data include the program SUPERMASSA, which models the relative intensity of the two
306 alternative alleles using Normal densities (Serang *et al.* 2012).

307 A second important factor for using our model is that, although estimates of allele frequencies can
308 be accurate when sequencing coverage is low and sample sizes are large, the distribution for genotypes
309 is likely going to be quite diffuse. For analyses that treat genotypes as a nuisance parameter, this is
310 not an issue since we can integrate across genotype uncertainty. However, if the genotype *is* of primary
311 interest, then the experimental design of the study will need to change to acquire higher coverage
312 at each locus for more accurate genotype estimation. Therefore, the decision between sequencing
313 more individuals with lower average coverage versus sequencing fewer individuals with higher average
314 coverage depends primarily on whether the genotypes will be used or not.

315 **Model adequacy**

316 As noted earlier, the probability model that we use for the inheritance of alleles is one of polysomy
317 without double reduction. In some cases, this model may be inappropriate but it can still be
318 informative to check for loci that do or do not follow the model that we assume. Below we describe a
319 simple procedure for rejecting our model of inheritance on a per locus basis using comparisons with
320 the posterior predictive distribution of sequencing reads. Model checking is an important part of

321 making statistical inferences and can play a role in understanding when a model adequately describes
322 the data being analyzed. In the case of our model, it can serve as a basis for understanding the
323 inheritance patterns of the organism being studied by determining which loci adhere to a simple
324 pattern of polysomic inheritance.

325 Given M posterior samples for the allele frequencies at locus ℓ , $\{p_\ell^{[1]}, p_\ell^{[2]}, \dots, p_\ell^{[M]}\}$, we will
326 simulate new values for the genotypes ($\tilde{g}_{i\ell}$) and reference read counts ($\tilde{r}_{i\ell}$) for all individuals and
327 use the ratio of simulated reference read counts to observed total read counts ($\frac{\tilde{r}_{i\ell}}{t_{i\ell}}$) as a summary
328 statistic for comparing the observed read count ratios to the distribution of the predicted read count
329 ratios. The use of the likelihood (or similar quantities) as a summary statistic has been a common
330 practice in posterior predictive comparisons of nucleotide substitution models, and more recently for
331 comparative phylogenetics (Ripplinger & Sullivan 2010; Reid *et al.* 2014; Pennell *et al.* 2015). We use
332 the ratio of reference to total read counts here because it is the maximum likelihood estimate of the
333 probability of success for a Binomial random variable and because it is a simple quantity to calculate.
334 The use of other summary statistics, or a combination of multiple summary statistics, would also
335 possible. The procedure for our posterior predictive model check is as follows:

- 336 1. For locus $\ell = 1, \dots, L$:
 - 337 1.1. For posterior sample $m = 1, \dots, M$:
 - 338 1.1.1. Simulate new genotype values ($\tilde{g}_{i\ell}^{[m]}$) for all individuals ($i = 1, \dots, N$) by drawing
339 from a Binomial ($\psi, p_\ell^{[m]}$).
 - 340 1.1.2. Simulate new reference read counts ($\tilde{r}_{i\ell}^{[m]}$) from each new genotype for all individuals
341 by drawing from Eq. 1.
 - 342 1.1.3. Calculate the reference read ratio for the simulated data for sample m and sum across
343 individuals: $\tilde{\mathcal{S}}_\ell^{[m]} = \sum_{i=1}^N \left(\frac{\tilde{r}_{i\ell}^{[m]}}{t_{i\ell}} \right)$.
 - 344 1.1.4. Calculate the reference read ratio for the observed data and sum across individuals:
345 $\mathcal{S}_\ell = \sum_{i=1}^N \left(\frac{r_{i\ell}}{t_{i\ell}} \right)$.
 - 346 1.2. Calculate the difference between the observed reference read ratio and the M simulated
347 reference read ratios: $\left\{ \mathcal{S}_\ell - \tilde{\mathcal{S}}_\ell^{[1]}, \dots, \mathcal{S}_\ell - \tilde{\mathcal{S}}_\ell^{[M]} \right\}$.
- 348 2. Determine if the 95% highest posterior density (HPD) interval of the distribution of re-centered
349 reference read ratios contains 0.

350 When the distribution of the differences in ratios between the observed and simulated datasets
351 does not contain 0 in the 95% HPD interval, it provides evidence that the locus being examined
352 does not follow a pattern of strict polysomic inheritance. A similar approach could be used on an
353 individual basis by comparing the observed ratio of reference reads to the predicted ratios for each
354 individual at each locus. We implement the per locus version of this posterior predictive model
355 checking procedure in the POLYFREQS package with the function `polyfreqs_pps`.

356 Extensibility

357 The modular nature of our hierarchical model allows for the addition and modification of nodes in the
358 model graph (Figure 1). One of the simplest extensions to the model that can build directly on the
359 current setup would be to consider loci with more than two alleles. This can be done using Multinomial
360 distributions for sequencing reads and genotypes and a Dirichlet prior on allele frequencies (the
361 Multinomial and Dirichlet distributions form a conjugate family; Gelman *et al.* 2014). We could also
362 model populations of mixed ploidy by using a vector of individually assigned ploidy levels instead of
363 assuming a single value for the whole population ($\psi = \{\psi_1, \dots, \psi_N\}$). However, this would assume
364 random mating among ploidy levels.

365 The place where we believe our model could have the greatest impact is through modifications
366 and extensions of the probability model used for the inheritance of alleles. These models have been
367 difficult to apply in the past as a result of genotype uncertainty. However, using our model as a
368 starting point, it could be possible to infer patterns of inheritance (polysomy, disomy, heterosomy)
369 and other demographic parameters (e.g., effective population size, population differentiation) without
370 requiring direct knowledge about the genotypes of the individuals in the population. For example,
371 Haldane's (1930) model of genotype frequencies for autopolyploids that are partially selfing could
372 be used to infer the prevalence of self-fertilization within a population. Similarly, Fisher's (1943)
373 model for double reduction in the inheritance of style lengths for *Lythrum* could be generalized and
374 used alone or together with a model for partial selfing to better understand how these processes
375 affect the genetic diversity of a population. A more recent model described by Stift *et al.* (2008)
376 used microsatellites to infer the different inheritance patterns (disomic, tetrasomic, intermediate) for
377 tetraploids in the genus *Rorippa* (Brassicaceae) following crossing experiments. The reformulation of
378 such a model for biallelic SNPs gathered using high throughput sequencing could provide a suitable
379 framework for understanding inheritance patterns across the genome. An ideal model would be one
380 that could help to understand inheritance patterns without the need conduct additional experiments.
381 However, to our knowledge, such a model does not currently exist and may not even be possible to
382 implement due to the complexity of possible inheritance patterns that might need to be considered
383 without the addition of information from crosses.

384 Conclusions

385 The recent emergence of models for genotype uncertainty in diploids has introduced a theoretical
386 framework for dealing with the fact that genotypes are unobserved quantities (Gompert & Buerkle
387 2012; Buerkle & Gompert 2013). Our extension of this theory to cases of higher ploidy (specifically
388 to autopolyploids) progresses naturally from the original work but also serves to alleviate the deeper
389 issue of ADU. The power and flexibility of these models as applied at the diploid level has the
390 potential to be replicated for polyploid organisms with the addition of suitable models for allelic
391 inheritance. The construction of hierarchical models containing suitable probability models for ADU,
392 allelic inheritance and perhaps even additional levels for important parameters such as F statistics
393 or the allele frequency spectrum also have the potential to provide key insights into the population
394 genetics of polyploids (Gompert & Buerkle 2011; Buerkle & Gompert 2013). Future work on such
395 models will help to progress the study of polyploid taxa and could eventually lead to more generalized
396 models for understanding the processes that have shaped their evolutionary histories.

397 Acknowledgements

398 The authors would like to thank the Ohio Supercomputer Center for access to computing resources
399 and Nick Skomrock for assistance with deriving the full conditional distributions of the model in
400 the diploid case. We would also like to thank *Associate Editor*, Aaron Wenzel and # anonymous
401 reviewers for their comments on the manuscript. This work was partially funded through a grant
402 from the National Science Foundation (DEB-1455399) to ADW and LSK.

403 References

- 404 Allendorf FW, Thorgaard GH (1984) *Tetraploidy and the evolution of salmonid fishes*. In: *Evolutionary*
405 *genetics of fishes*. Edited by B. J. Turner. Plenum Press, pp. 1–53.
- 406 Arnold B, Bomblies K, Wakeley J (2012) Extending coalescent theory to autotetraploids. *Genetics*,
407 **192**, 195–204.
- 408 Baird NA, Etter PD, Atwood TS, *et al.* (2008) Rapid SNP discovery and genetic mapping using
409 sequenced RAD markers. *PLoS ONE*, **3**, e3376.
- 410 Barlow N (1913) Preliminary note on heterostylism in *Oxalis* and *Lythrum*. *Journal of Genetics*, **3**,
411 53–65.
- 412 Barlow N (1923) Inheritance of the three forms in trimorphic plants. *Journal of Genetics*, **13**,
413 133–146.
- 414 Buerkle CA, Gompert Z (2013) Population genomics based on low coverage sequencing: how low
415 should we go? *Molecular Ecology*, **22**, 3028–3035.
- 416 Cannon SB, McKain MR, Harkess A, *et al.* (2014) Multiple polyploidy events in the early radiation
417 of nodulating and nonnodulating legumes. *Molecular Biology and Evolution*, **32**, 193–210.
- 418 Clark LV, Jasieniuk M (2011) POLYSAT: an R package for polyploid microsatellite analysis. *Molecular*
419 *Ecology Resources*, **11**, 562–566.
- 420 Clausen J, Keck DD, Hiesey WM (1940) *Experimental studies on the nature of species. I. Effect of*
421 *varied environments on western American plants*. Carnegie Inst. Washington Publ.
- 422 Clausen J, Keck DD, Hiesey WM (1945) *Experimental studies on the nature of species. II. Plant*
423 *evolution through amphiploidy and autopolyploidy, with examples from Madiinae*. Carnegie Inst.
424 Washington Publ.
- 425 Cui L, Wall PK, Leebens-Mack JH, *et al.* (2006) Widespread genome duplications throughout the
426 history of flowering plants. *Genome Research*, **16**, 738–749.
- 427 Dufresne F, Stift M, Vergilino R, Malbe BK (2014) Recent progress and challenges in population
428 genetics of polyploid organisms: an overview of current state-of-the-art molecular and statistical
429 tools. *Molecular Ecology*, **23**, 40–69.
- 430 East EM (1927) The inheritance of heterostyly in *Lythrum salicaria*. *Genetics*, **12**, 393–414.

- 431 Eddelbuettel D (2013) *Seamless R and C++ integration with Rcpp*. Springer, New York.
- 432 Eddelbuettel D, François R (2011) Rcpp: seamless R and C++ integration. *Journal of Statistical*
433 *Software*, **40**, 1–18.
- 434 Eddelbuettel D, Sanderson C (2014) RcppArmadillo: accelerating R with high-performance C++
435 linear algebra. *Computational Statistics and Data Analysis*, **71**, 1054–1063.
- 436 Esselink GD, Nybom H, Vosman B (2004) Assignment of allelic configuration in polyploids using
437 the MAC-PR (microsatellite DNA allele counting–peak ratios) method. *Theoretical and Applied*
438 *Genetics*, **109**, 402–408.
- 439 Fisher RA (1941) The theoretical consequences of polyploid inheritance for the Mid style form of
440 *Lythrum salicaria*. *Annals of Eugenics*, **11**, 31–38.
- 441 Fisher RA (1943) Allowance for double reduction in the calculation of genotype frequencies with
442 polysomic inheritance. *Annals of Eugenics*, **12**, 169–171.
- 443 Fisher RA, Mather K (1943) The inheritance of style length in *Lythrum salicaria*. *Annals of Eugenics*,
444 **12**, 1–23.
- 445 Furlong RF, Holland PWH (2001) Were vertebrates octoploid? *Philosophical Transactions of the*
446 *Royal Society B: Biological Sciences*, **357**, 531–544.
- 447 Gelman A, Carlin JB, Stern HS, Dunson DB, Vehtari A, Rubin DB (2014) *Bayesian data analysis*.
448 Chapman & Hall/CRC Press, 3rd edn.
- 449 Gompert Z, Buerkle CA (2011) A hierarchical Bayesian model for next-generation population
450 genomics. *Genetics*, **187**, 903–917.
- 451 Gompert Z, Buerkle CA (2012) bgc: software for Bayesian estimation of genomic clines. *Molecular*
452 *Ecology Resources*, **12**, 1168–1176.
- 453 Grant V (1971) *Plant speciation*. Columbia University Press.
- 454 Gregory TR, Mable BK (2005) *Polyploidy in animals*. In: *The evolution of the genome*. Edited by T.
455 R. Gregory. Elsevier, pp. 427–517.
- 456 Haldane JBS (1930) Theoretical genetics of autopolyploids. *Journal of Genetics*, **22**, 359–372.
- 457 Jiao Y, Wickett NJ, Ayyampalayam S, *et al.* (2011) Ancestral polyploidy in seed plants and
458 angiosperms. *Nature*, **473**, 97–100.
- 459 Jones G, Sagitov S, Oxelman B (2013) Statistical inference of allopolyploid species networks in the
460 presence of incomplete lineage sorting. *Systematic Biology*, **62**, 467–478.
- 461 Kimura M (1964) Diffusion models in population genetics. *Journal of Applied Probability*, **1**, 177–232.
- 462 Lawrence WJC (1929) The genetics and cytology of *Dahlia* species. *Journal of Genetics*, **21**, 125–158.
- 463 Logan-Young CJ, Yu JZ, Verma SK, Percy RG, Pepper AE (2015) SNP discovery in complex
464 allotetraploid genomes (*Gossypium* spp., Malvaceae) using genotyping by sequencing. *Applications*
465 *in Plant Sciences*, **3**, 1400077.

- 466 Miller MR, Dunham JP, Amores A, Cresko WA, Johnson EA (2007) Rapid and cost-effective
467 polymorphism identification and genotyping using restriction site associated DNA (RAD) markers.
468 *Genome Research*, **17**, 240–248.
- 469 Moody ML, Mueller LD, Soltis DE (1993) Genetic variation and random drift in autotetraploid
470 populations. *Genetics*, **134**, 649–657.
- 471 Muller HJ (1914) A new mode of segregation in gregory’s tetraploid primulas. *American Naturalist*,
472 **48**, 508–512.
- 473 Ogden R, Gharbi K, Mugue N, *et al.* (2013) Sturgeon conservation genomics: SNP discovery and
474 validation using RAD sequencing. *Molecular Ecology*, **22**, 3112–3123.
- 475 Ohno S (1970) *Evolution by gene duplication*. Springer.
- 476 Otto SP, Whitton J (2000) Polyploid incidence and evolution. *Annual Review of Genetics*, **34**,
477 401–437.
- 478 Parisod C, Holderegger R, Brochmann C (2010) Evolutionary consequences of autopolyploidy. *New*
479 *Phytologist*, **186**, 5–17.
- 480 Pennell MW, FitzJohn RG, Cornwell WK, Harmon LJ (2015) Model adequacy and the macroevolution
481 of angiosperm functional traits. *American Naturalist*, **186**, E100.
- 482 Plummer M, Best N, Cowles K, Vines K (2006) CODA: Convergence Diagnostics and Output Analysis
483 for MCMC. *R News*, **6**, 7–11.
- 484 Puritz JB, Matz MV, Toonen RJ, Weber JN, Bolnick DI, Bird CE (2014) Demystifying the RAD
485 fad. *Molecular Ecology*, **23**, 5937–5942.
- 486 R Core Team (2014) *R: a language and environment for statistical computing*. R Foundation for
487 Statistical Computing, Vienna, Austria.
- 488 Ramsey J (2011) Polyploidy and ecological adaptation in wild yarrow. *Proceedings of the National*
489 *Academy of Sciences*, **108**, 7096–7101.
- 490 Ramsey J, Ramsey TS (2014) Ecological studies of polyploidy in the 100 years following its discovery.
491 *Philosophical Transactions of the Royal Society B: Biological Sciences*, **369**, 20130352.
- 492 Reid NM, Hird SM, Brown JM, *et al.* (2014) Poor fit to the multispecies coalescent is widely detectable
493 in empirical data. *Systematic Biology*, **63**, 322–333.
- 494 Ripplinger J, Sullivan J (2010) Assessment of substitution model adequacy using frequentist and
495 Bayesian methods. *Molecular Biology and Evolution*, **27**, 2790–2803.
- 496 Scarpino SV, Levin DA, Meyers LA (2014) Polyploid formation shapes flowering plant diversity.
497 *American Naturalist*, **184**, doi: 10.1086/677752.
- 498 Selmecki AM, Maruvka YE, Richmond PA, *et al.* (2015) Polyploidy can drive rapid adaptation in
499 yeast. *Nature*, **519**, 349–352.

- 500 Serang O, Mollinari M, Garcia AAF (2012) Efficient exact maximum *a posteriori* computation for
501 Bayesian SNP genotyping in polyploids. *PLoS ONE*, **7**, e30906.
- 502 de Silva H, Hall A, Rikkerink E, McNeilage M, Fraser L (2005) Estimation of allele frequencies in
503 polyploids under certain patterns of inheritance. *Heredity*, **95**, 327–334.
- 504 Soltis DE, Albert VA, Leebens-Mack J, *et al.* (2009) Polyploidy and angiosperm diversification.
505 *American Journal of Botany*, **96**, 336–348.
- 506 Soltis DE, Buggs RJA, Doyle JJ, Soltis PS (2010) What we still don't know about polyploidy. *Taxon*,
507 **59**, 1387–1403.
- 508 Soltis DE, Soltis PS, Schemske DW, *et al.* (2007) Autopolyploidy in angiosperms: have we grossly
509 underestimated the number of species? *Taxon*, **56**, 13–30.
- 510 Soltis DE, Soltis PS, Tate JA (2003) Advances in the study of polyploidy since plant speciation. *New*
511 *Phytologist*, **161**, 173–191.
- 512 Soltis DE, Visger CJ, Soltis PS (2014) The polyploidy revolution then...and now: Stebbins revisited.
513 *American Journal of Botany*, **101**, 1057–1078.
- 514 Soltis PS, Soltis DE (2009) The role of hybridization in plant speciation. *Annual Review of Plant*
515 *Biology*, **60**, 561–588.
- 516 Stebbins GL (1950) *Variation and evolution in plants*. Columbia University Press.
- 517 Stift M, Berenos C, Kuperus P, van Tienderen PH (2008) Segregation models for disomic, tetrasomic
518 and intermediate inheritance in tetraploids: a general procedure applied to *Rorippa* (yellow cress)
519 microsatellite data. *Genetics*, **179**, 2113–2123.
- 520 Wagner WH (1970) Biosystematics and evolutionary noise. *Taxon*, **19**, 146–151.
- 521 Wang N, Thomson M, Bodles WJA, *et al.* (2013) Genome sequence of dwarf birch (*Betula nana*)
522 and cross-species RAD markers. *Molecular Ecology*, **22**, 3098–3111.
- 523 Wickham H (2007) Reshaping data with the reshape package. *Journal of Statistical Software*, **21**,
524 1–20.
- 525 Wickham H (2009) *ggplot2: elegant graphics for data analysis*. Springer, New York.
- 526 Wickham H (2011) The split-apply-combine strategy for data analysis. *Journal of Statistical Software*,
527 **40**, 1–29.
- 528 Winge Ö (1917) The chromosomes: their number and general importance. *Compt. Rend. Trav. Lab.*
529 *Carlsberg*, **13**, 131–275.
- 530 Winkler H (1916) Über die experimentelle Erzeugung von Pflanzen mit abweichenden
531 Chromosomenzahlen. *Zeitschr. f. Bot.*, **8**, 417–531.
- 532 Wood TE, Takebayashi N, Barker MS, Mayrose I, Greenspoon PB, Rieseberg LH (2009) The frequency
533 of polyploid speciation in vascular plants. *Proceedings of the National Academy of Sciences*, **106**,
534 13875–13879.

535 Wright S (1931) Evolution in Mendelian populations. *Genetics*, **16**, 97–159.

536 Wright S (1937) The distribution of gene frequencies in populations. *Proceedings of the National*
537 *Academy of Sciences*, **23**, 307–320.

538 Wright S (1938) The distribution of gene frequencies in populations of polyploids. *Proceedings of the*
539 *National Academy of Sciences*, **24**, 372–377.

540 **Author Contributions**

541 Conceived of the study: PDB, LSK and ADW. PDB derived the polyploid model, ran the simulations,
542 coded the R package and wrote the initial draft of the manuscript. PDB, LSK and ADW reviewed
543 all parts of the manuscript and all authors approved of the final version.

544 **Data Accessibility**

545 Scripts for simulating the datasets, analyzing them using Gibbs sampling and producing the
546 figures from the resulting output can all be found on GitHub ([https://github.com/pblischak/](https://github.com/pblischak/polyfreqs-ms-data)
547 [polyfreqs-ms-data](https://github.com/pblischak/polyfreqs-ms-data)). We also provide an implementation of the Gibbs sampler for estimating
548 allele frequencies in the developmental R package POLYFREQS ([https://github.com/pblischak/](https://github.com/pblischak/polyfreqs)
549 [polyfreqs](https://github.com/pblischak/polyfreqs)). See the package vignette or GitHub wiki for more details ([https://github.com/](https://github.com/pblischak/polyfreqs/wiki)
550 [pblischak/polyfreqs/wiki](https://github.com/pblischak/polyfreqs/wiki)).

Table 1: Notation and symbols used in the description of the model for estimating allele frequencies in polyploids. Vector and matrix forms of the variables are also provided when appropriate.

Symbol	Description
L	The number of loci.
ℓ	Index for loci ($\ell \in \{1, \dots, L\}$).
N	Total number of individuals sequenced.
i	Index for individuals ($i \in \{1, \dots, N\}$).
ψ	The ploidy level of individuals in the population (e.g., tetraploid: $\psi=4$).
p_ℓ	Frequency of the reference allele at locus ℓ . [\mathbf{p}]
$g_{i\ell}$	The number of copies of the reference allele for individual i at locus ℓ . [\mathbf{G}]
$\tilde{g}_{i\ell}$	Simulated genotype for posterior predictive model checking.
$t_{i\ell}$	The total number of reads for individual i at locus ℓ . [\mathbf{T}]
$r_{i\ell}$	The number of reads with the reference allele for individual i at locus ℓ . [\mathbf{R}]
$\tilde{r}_{i\ell}$	Simulated reference read count for posterior predictive model checking. [$\tilde{\mathbf{R}}$]
ϵ	Sequencing error.

Figure 1: Graphical representation of a hierarchical Bayesian model for estimating allele frequencies. The two levels (allelic dosage and inheritance) represent the probability models that are used for inference from one graph to another. The model we present here focuses primarily on allelic dosage.

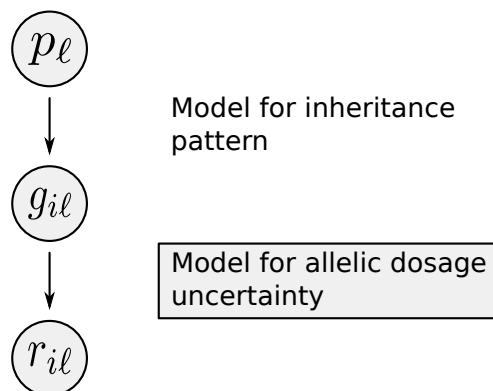


Figure 2: Heat maps of error rates (standard deviation from the true value) for allele frequency estimation in tetraploids. The x-axis shows the number of individuals (i5, i10, i20, i30) and the y-axis represents the sequencing coverage (c5, c10, c20, c50, c100) for each simulation. Note that the scales for each heat map are not the same, but the overall pattern of increased accuracy as the number of individuals and sequencing coverage increases for all allele frequencies.

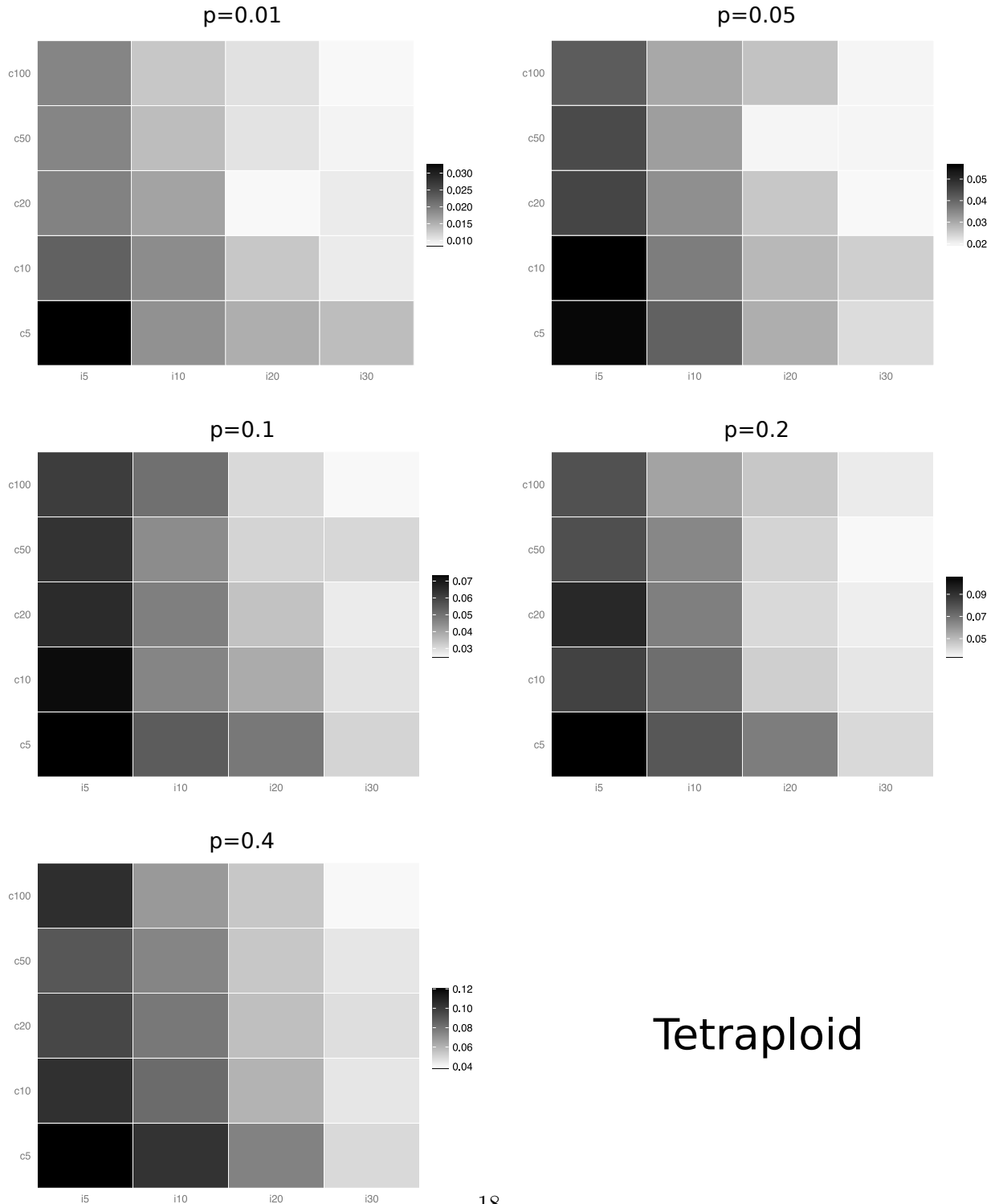


Figure 3: Heat maps of error rates (standard deviation from the true value) for allele frequency estimation in hexaploids. The x-axis shows the number of individuals (i5, i10, i20, i30) and the y-axis represents the sequencing coverage (c5, c10, c20, c50, c100) for each simulation. Note that the scales for each heat map are not the same, but the overall pattern of increased accuracy as the number of individuals and sequencing coverage increases is the same for all allele frequencies.

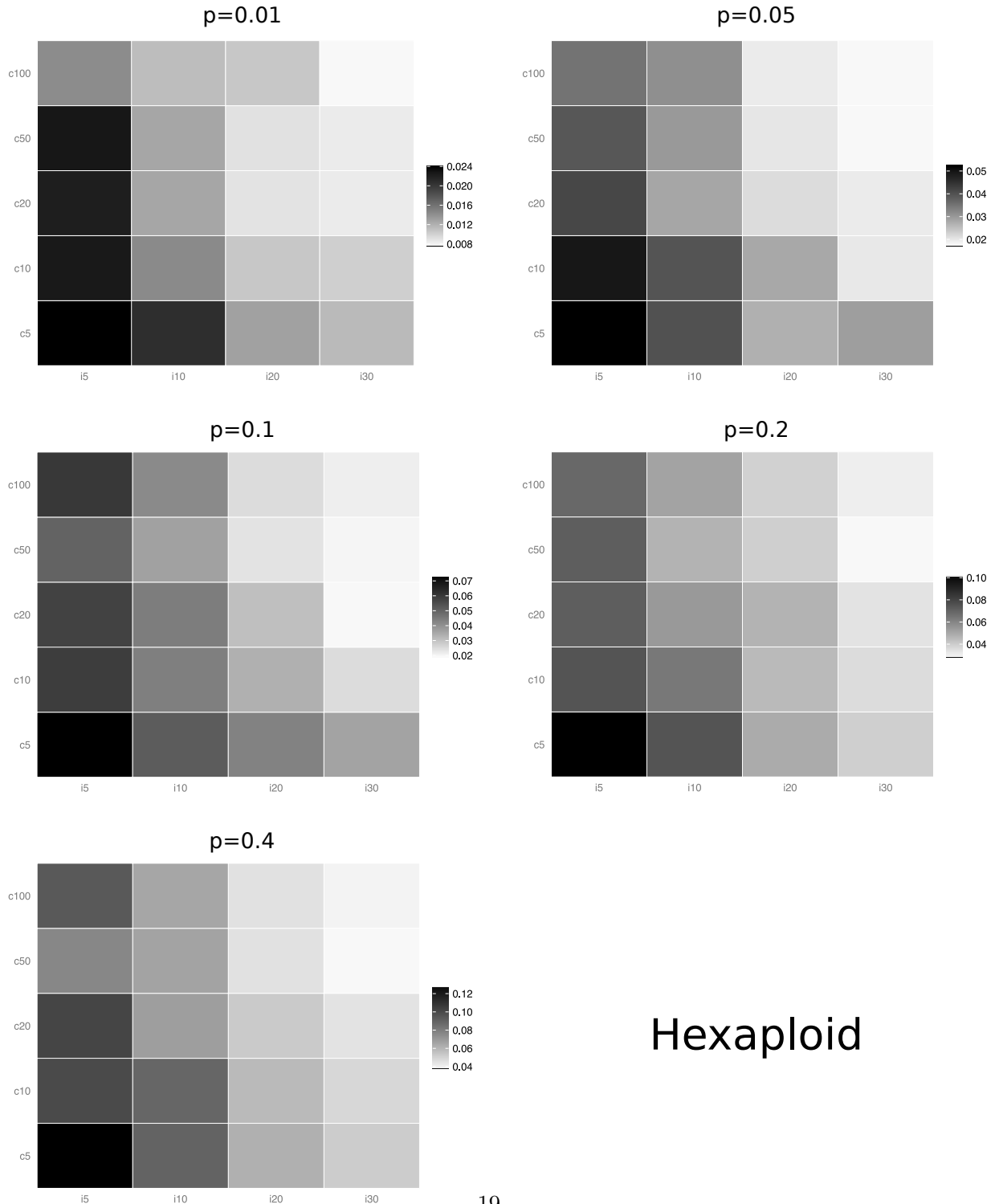


Figure 4: The posterior standard deviation for allele frequencies decreases with increased sequencing coverage. This plot provides a comparison of the distribution of posterior standard deviations of the 100 replicates performed for each level of sequencing coverage for the hexaploid simulation with 30 individuals and an allele frequency of 0.2.

