

# Performing parentage analysis for polysomic inheritances based on allelic phenotypes

Kang Huang,<sup>1,2</sup> Gwendolyn Huber,<sup>2</sup> Kermit Ritland,<sup>2</sup> Derek W. Dunn,<sup>1</sup> and Baoguo Li<sup>1,3,\*</sup>

<sup>1</sup>Shaanxi Key Laboratory for Animal Conservation, College of Life Sciences, Northwest University, Xi'an 710069, China

<sup>2</sup>Department of Forest and Conservation Sciences, University of British Columbia, Vancouver, BC V6T1Z4, Canada

<sup>3</sup>Center for Excellence in Animal Evolution and Genetics, Chinese Academy of Sciences, Kunming 650223, China

\*Corresponding author: Shaanxi Key Laboratory for Animal Conservation, College of Life Sciences, Northwest University, Xi'an 710069, China. baoguoli@nwu.edu.cn

## Abstract

Polyploidy poses several problems for parentage analysis. We present a new polysomic inheritance model for parentage analysis based on genotypes or allelic phenotypes to solve these problems. The effects of five factors are simultaneously accommodated in this model: (1) double-reduction, (2) null alleles, (3) negative amplification, (4) genotyping errors and (5) self-fertilization. To solve genotyping ambiguity (unknown allele dosage), we developed a new method to establish the likelihood formulas for allelic phenotype data and to simultaneously include the effects of our five chosen factors. We then evaluated and compared the performance of our new method with three established methods by using both simulated data and empirical data from the cultivated blueberry (*Vaccinium corymbosum*). We also developed and compared the performance of two additional estimators to estimate the genotyping error rate and the sample rate. We make our new methods freely available in the software package POLYGENE, at <http://github.com/huangkang1987/polygene>.

**Keywords:** parentage analysis; polysomic inheritance; genotyping ambiguity; double-reduction; null alleles; self-fertilization

## Introduction

Parentage analysis is a common technique in plant ecology and selective breeding. This technique for identifying parents enables researchers to assess seed dispersal (Ismail *et al.* 2017), pollen dispersal (Bezemer *et al.* 2016), assortative mating (Monthe *et al.* 2017), isolation (Tambarussi *et al.* 2015), current gene flow (Duminil *et al.* 2016), mating systems (Tan *et al.* 2019), reproductive success (Watanabe *et al.* 2018), functional sex (Oddou-Muratorio *et al.* 2018), and to increase genetic gain from selective breeding (Norman *et al.* 2018).

A large proportion of plant species are polyploid, with 24% of all plant taxa showing some form of polysomic inheritance (Barker *et al.* 2016), and at least 47% of angiosperm species having polyploidy in their ancestral lineage (Wood *et al.* 2009). Existing methods of parentage analysis for polyploids use the pseudo-dominant approach (Rodzen *et al.* 2004; Wang and Scribner 2014) and exclusion approach (Zwart *et al.* 2016). In the pseudo-dominant approach, the polyploid genotypes or the allelic phenotypes are converted into pseudo-dominant phenotypes and use diploid likelihood equations to calculate the likelihood for parentage assignment (Gerber *et al.* 2000), in which each allele at a codominant locus is treated as an independent dominant "locus." This approach enables rapid calculation but is inferior to that based on polysomic inheritance methods because any transformation of data will cause a loss of information and thus a reduction in accuracy (Wang and Scribner 2014). The exclusion approach excludes the parents based on Mendelian

incompatibility. However, due to the high gamete diversity (Pelé *et al.* 2018) and genotyping ambiguity (Huang *et al.* 2014), the exclusion rate is low in polyploid, especially for a parent-offspring pair. Thus, the development of more accurate methods of parentage analysis for polyploids is required.

Several models for polysomic inheritance have been developed, such as double-reduction models (Muller 1914; Haldane 1930; Mather 1935), genotypic frequencies (Fisher 1943; Geiringer 1949), and transitional probabilities from a zygote to a gamete (Fisher and Mather 1943; Field *et al.* 2017). On the basis of these findings, Huang *et al.* (2019) derived the generalized genotypic frequency and gamete frequency for ploidy levels fewer than 12 and derived the generalized transitional probability from a zygote to a gamete for any ploidy level. These models provide a foundation on which to establish a method of parentage analysis for polyploids.

A unique feature of polysomic inheritance is *double-reduction* such that a pair of sister chromatids are segregated into a single gamete (Parisod *et al.* 2010). Double-reduction arises from a combination of three major events during meiosis: (1) the crossing-over between non-sister chromatids, (2) an appropriate pattern of disjunction, and (3) the migration of chromosomal segments carrying a pair of sister chromatids to the same gamete (Darlington 1929; Haldane 1930). Geneticists have developed several mathematical models to simulate double-reduction: these are the *random chromosome segregation* (RCS) model (i.e. without double-reduction; Muller 1914), the *pure random chromatid segregation* (PRCS) model (Haldane 1930), the *complete equational*

Received: September 13, 2020. Accepted: November 09, 2020

© The Author(s) 2021. Published by Oxford University Press on behalf of Genetics Society of America.

This is an Open Access article distributed under the terms of the Creative Commons Attribution-NonCommercial-NoDerivs licence (<http://creativecommons.org/licenses/by-nc-nd/4.0/>), which permits non-commercial reproduction and distribution of the work, in any medium, provided the original work is not altered or transformed in any way, and that the work is properly cited. For commercial re-use, please contact [journals.permissions@oup.com](mailto:journals.permissions@oup.com)

segregation (CES) model (Mather 1935) and the *partial equational segregation* (PES) model (Huang et al. 2019). A brief description of each of these models is given in Supplementary Appendix A.

There are two consequences of double-reduction that will influence parentage analysis: (1) the genotypic frequencies will deviate from expected values under RCS, resulting in a bias of the estimated LOD scores and (2) some unexpected offspring genotypes may be generated (e.g. an offspring genotype AAEE is produced from ABCD × EFGH) along with the true father being excluded. Therefore, the complete array of diverse polyploid offspring genotypes has to be accounted for in order to conduct a comprehensive and accurate paternity analysis (Stift et al. 2008, 2010).

There are also several additional problems associated with PCR-based markers that need to be accounted, irrespective of ploidy. One problem is the *genotyping ambiguity* of polyploids (Huang et al. 2014), in the sense that the allelic dosage of PCR-based markers cannot be determined. For example, the genotype AABB will appear to be identical to AAAB. Another problem arises when using microsatellites, which are the genetic markers most frequently used for parentage analysis. Microsatellites can have null alleles (Ravinet et al. 2016) that cause both the lack of amplification of null allele homozygotes and the lack of detectability of null allele heterozygotes (Wagner et al. 2006). A third problem comes from genotyping errors, which may cause a true parent to be mistakenly excluded due to an observed lack of shared alleles with the offspring (Blouin 2003). Finally, inbreeding will result in an excess of homozygotes in a population, such as when plants self-fertilize (Ritland 2002). The genotypic frequencies used for a parentage analysis will thus be affected by any inbreeding.

Here, we extend the disomic inheritance model of Kalinowski et al. (2007) to account for polysomic inheritance to enable accurate parentage analysis for polyploids based on genotypes or allelic phenotypes. Our new polysomic inheritance model accommodates the effects of five factors: (1) double-reduction, (2) null alleles, (3) negative amplification, (4) genotyping errors and (5) self-fertilization. To solve the problem of genotyping ambiguity, we develop a new method so as to establish the likelihood formulas for allelic phenotype data, with the effects of our five factors of interest also being included in these formulas. We subsequently use a designated simulated dataset to evaluate and compare the performance of our new method with three other established methods. We also use an empirical microsatellite dataset from the cultivated blueberry (*Vaccinium corymbosum*) to test the performance of all four methods. Moreover, we develop and evaluate two models to estimate the genotyping error rate and the sample rate (the probability that a true parent is sampled). We have incorporated our new parentage analysis methods into the software package POLYGENE, which can be freely downloaded at <http://github.com/huangkang1987/polygene>.

## Methods

Here we assume that our parentage analysis model satisfies four assumptions, which are also commonly used for diploid population genetics methods. These four assumptions are: (1) the population is large enough to neglect any effects of genetic drift and there is no population subdivision; (2) that mating is random, except for a given rate of selfing, and the probability of selfing is the same among individuals (3) the distributions of the genotypes are the same for males and females, and reach an equilibrium state (i.e. genotypic frequencies do not change among generations), and (4) the genetic markers used are neutral, autosomal, codominant and unlinked.

The multiset consisting of allele copies within an individual at a locus is called a *genotype*, denoted by  $\mathcal{G}$  or  $G$ , in which  $\mathcal{G}$  represents an observed genotype and  $G$  represents a true genotype. For example, {A, A, A, B} is a genotype, abbreviated as AAAB. The set consisting of alleles within an individual at a locus is called an *allelic phenotype*, or a *phenotype* for short, denoted by  $\mathcal{P}$ . For instance, {A, B} is a phenotype, written as AB for short.

Our methods are the extensions of Kalinowski et al.'s (2007) method. In the following text, we briefly describe the scheme of Kalinowski et al.'s (2007) method and its associated diploid model.

## Scheme of simulation-based likelihood approach

The foundations for assigning parentage with confidence by a simulation-based likelihood approach were established by Marshall et al. (1998). There are three typical categories in this approach: (1) identifying the father (or one parent) when the mother (or the other parent) is unknown; (2) identifying the father (or one parent) when the mother (or the other parent) is known; and (3) identifying the father and the mother (or parents) jointly. There are two situations in the third category, the first is for dioecious species and the sexes of individuals are recorded (termed sexes known), and the second is for monoecious species or the sexes of individuals are not recorded (termed sexes unknown). The procedures of a parentage analysis are broadly as follows.

For each of the first two categories, two hypotheses are established: the *first hypothesis* is that the alleged father is the true father, denoted by  $H_1$ ; the *alternative hypothesis* is that the alleged father is not the true father, denoted by  $H_2$ . For the third category, "father" needs to be changed to "parents" in both hypotheses.

Given a hypothesis  $H$ , the *likelihood* is defined as the probability of some observed data given  $H$ , written as  $\mathcal{L}(H)$ . Returning to  $H_1$  and  $H_2$  as described above, we call the natural logarithm of the ratio of  $\mathcal{L}(H_1)$  to  $\mathcal{L}(H_2)$  the *LOD score*, or *LOD* as the abbreviation, symbolically  $\text{LOD} = \ln \frac{\mathcal{L}(H_1)}{\mathcal{L}(H_2)}$ . Moreover, if a LOD is positive, it means that  $H_1$  is more likely to be true than  $H_2$ . Similarly, a negative LOD means that  $H_2$  is more likely to be true than  $H_1$ .

Marshall et al. (1998) provided a statistic  $\Delta$  for resolving paternity, the definition of which is:

$$\Delta = \begin{cases} \text{LOD}_1 - \text{LOD}_2 & \text{if } n \geq 2, \\ \text{LOD}_1 & \text{if } n = 1, \\ \text{undefined} & \text{if } n = 0, \end{cases}$$

where  $\text{LOD}_1$  and  $\text{LOD}_2$  are, respectively, the LODs of the most-likely and the next most-likely alleged fathers, and  $n$  is the number of all alleged fathers. For a practical application, the statistic  $\Delta$  needs to be singly calculated for each individual offspring. Monte-Carlo simulations are subsequently used to assess the confidence level of  $\Delta$ . The symbol  $\Delta_{0.95}$  represents that the threshold of  $\Delta$  reaches the confidence level 95%, in the sense that if  $\Delta \geq \Delta_{0.95}$ , the probability that the assigned parent is the true parent is at least 0.95.

The likelihood equations used in Marshall et al. (1998) to accommodate genotyping error miscalculate the probability of observing an erroneous genotype. Therefore, we applied the corrected equations in Kalinowski et al. (2007) in the following.

## Marshall et al.'s (1998) diploid model

Marshall et al.'s (1998) diploid model (abbreviated as the Ma-model) accounts for any genotyping errors under the assumption that the genotype frequencies accord with the Hardy-Weinberg equilibrium (HWE). This model consists of some likelihood formulas (listed in the first half of Supplementary Appendix B)

together with the rules and methods for a general parentage analysis.

The likelihood formulas of the Ma-model are derived by using the transitional probability  $T(\mathcal{G} | G)$  from a true genotype  $G$  to an observed genotype  $\mathcal{G}$ , whose expression is:

$$T(\mathcal{G} | G) = (1 - e)\mathcal{B}_{\mathcal{G}=G} + e\Pr(\mathcal{G}), \quad (1)$$

where  $e$  is the genotyping error rate,  $\Pr(\mathcal{G})$  is the frequency of  $\mathcal{G}$ , and  $\mathcal{B}_X$  is a binary variable, such that  $\mathcal{B}_X = 1$  if the expression  $X$  is true, or  $\mathcal{B}_X = 0$  otherwise.

As previously stated, the procedures underlying the Ma-model to perform a parentage analysis are as follows: (1) calculating  $\mathcal{L}(H_1)$  and  $\mathcal{L}(H_2)$ , (2) finding the threshold of  $\Delta$ , (3) calculating the LOD and  $\Delta$ , and (4) using the values obtained in the previous three steps to assess the confidence level of this parentage analysis.

In the following text, we will use the first category in a parentage analysis as an example to show how to calculate the likelihoods  $\mathcal{L}(H_1)$  and  $\mathcal{L}(H_2)$  in the Ma-model. The expressions of  $\mathcal{L}(H_1)$  and  $\mathcal{L}(H_2)$  are:

$$\begin{aligned} \mathcal{L}(H_1) &= \Pr(\mathcal{G}_A)[(1 - e)^2 T(\mathcal{G}_O | \mathcal{G}_A) + 2e(1 - e)\Pr(\mathcal{G}_O) + e^2\Pr(\mathcal{G}_O)], \\ \mathcal{L}(H_2) &= \Pr(\mathcal{G}_A)\Pr(\mathcal{G}_O), \end{aligned} \quad (2)$$

where  $\mathcal{G}_A$  and  $\mathcal{G}_O$  are, respectively, the observed genotypes of the alleged father and the offspring,  $\Pr(\mathcal{G}_A)$  and  $\Pr(\mathcal{G}_O)$  are their frequencies, and  $T(\mathcal{G}_O | \mathcal{G}_A)$  is the transitional probability from  $\mathcal{G}_A$  to  $\mathcal{G}_O$ .

In the Ma-model, the genotyping error is considered as the replacement of a true genotype with a random genotype according to the genotypic frequencies. Thus the genotyping error does not change the distribution of the genotypes, i.e.  $\Pr(\mathcal{G}) = \Pr(G = \mathcal{G})$ . Moreover,  $\Pr(G)$  can be directly calculated from the HWE prediction:

$$\Pr(G) = \begin{cases} p_i^2 & \text{if } G = A_i A_i, \\ 2p_i p_j & \text{if } G = A_i A_j. \end{cases}$$

This is because any null alleles, any negative amplification (i.e. amplification failure due to experimental error or a poor DNA quality, rather than a null allelic homozygote) and any inbreeding/selfing are not considered in the Ma-model.

Next, the transitional probability  $T(\mathcal{G}_O | \mathcal{G}_A)$  is calculated under the assumptions that  $\mathcal{G}_A$  and  $\mathcal{G}_O$  are correctly typed and that the alleged father is the true father, i.e. under the assumptions that  $G_O = \mathcal{G}_O$  and  $G_F = \mathcal{G}_A$ , where  $G_O$  and  $G_F$  are the true genotypes of the offspring and the true father, respectively. Therefore,  $T(\mathcal{G}_O | \mathcal{G}_A)$  is the same as  $T(G_O | G_F)$  under these assumptions. Because one allele within  $G_O$  is randomly inherited from the parents, and the other is randomly sampled from the population according to the allele frequencies, the transitional probability  $T(G_O | G_F)$  can be expressed as:

$$T(G_O | G_F) = \begin{cases} p_i & \text{if } G_O = A_i A_i \text{ and } G_F = A_i A_i, \\ p_j & \text{if } G_O = A_i A_j \text{ and } G_F = A_i A_i, \\ \frac{1}{2}(p_i + p_j) & \text{if } G_O = A_i A_j \text{ and } G_F = A_i A_j, \\ \frac{1}{2}p_k & \text{if } G_O = A_i A_k \text{ and } G_F = A_i A_j, \\ 0 & \text{otherwise,} \end{cases}$$

where  $A_i$ ,  $A_j$  and  $A_k$  are distinct identical-by-state alleles,  $p_i$ ,  $p_j$  and  $p_k$  are their frequencies.

Now, we see that the two likelihood formulas in Equation (2) can be used for the actual calculation as long as the values of the genotyping error rate  $e$  and those frequencies of alleles are given.

For the second and third categories in a parentage analysis, to calculate the transitional probabilities  $T(\mathcal{G}_O | \mathcal{G}_A, \mathcal{G}_M)$  and  $T(\mathcal{G}_O | \mathcal{G}_A, \mathcal{G}_{AM})$  in the likelihood formulas in the Ma-model (see the first half of Supplementary Appendix B), we need to apply the transitional probability  $T(G_O | G_F, G_M)$  from a pair of true genotypes of the true parents to a true genotype of the offspring. Because the genotypic frequencies in the Ma-model are in HWE, according to the Mendelian segregation (i.e. each parent randomly contributes one allele to an offspring genotype),  $T(\mathcal{G}_O | \mathcal{G}_A, \mathcal{G}_M)$  can be calculated by:

$$T(G_O | G_F, G_M) = \frac{1}{4} \sum_{i=1}^2 \sum_{j=1}^2 \mathcal{B}_{G_O=A_i B_j},$$

where  $A_i$  (or  $B_j$ ) is an allele within  $G_F$  (or  $G_M$ ).

## Polyplod model

The polysomic inheritance model (abbreviated as the *polyplod model*) presented here is for use with even levels of ploidy, and consists of some likelihood formulas and some additional conditions along with the rules and methods for a general parentage analysis. These additional conditions are: (1) which of the two data types (genotypic and phenotypic) are to be selected, (2) whether self-fertilization is considered, (3) whether null alleles and/or negative amplifications are to be considered, and (iv) which of the four double-reduction models, listed in Supplementary Table S1, is chosen.

As for the Ma-model, our new model accommodates the effect of genotyping errors and the presence of these errors will not change the genotypic and phenotypic frequencies. Moreover, if self-fertilization is considered in our model, its effect will also be incorporated into the likelihood formulas.

For the genotypic data, the likelihood formulas for all three categories in a parentage analysis, under either self-fertilization or not, are given in Supplementary Appendix B. For polysomic inheritance, the genotypic frequencies ( $\Pr(\mathcal{G})$ ) and transitional probabilities ( $T(G_O | G_F)$  and  $T(G_O | G_F, G_M)$ ) need to be properly adjusted, where the formula of  $\Pr(\mathcal{G})$  under inbreeding and double-reduction is given in Supplementary Appendix C [or in Huang et al. (2019)], and the formulas of  $T(G_O | G_F)$  and  $T(G_O | G_F, G_M)$  are given in Supplementary Appendix D.

For the phenotypic data, the likelihood formulas for all three categories in a parentage analysis under the condition of either self-fertilization or not are given in Supplementary Appendix E. In such circumstances, the phenotypic frequencies ( $\Pr(\mathcal{P})$ ) in these formulas are calculated by Equation (A5), and the transitional probabilities ( $T(\mathcal{P}_O | \mathcal{P}_F)$  and  $T(\mathcal{P}_O | \mathcal{P}_F, \mathcal{P}_M)$ ) by Equation (3) or (4). To solve the problem of genotyping ambiguity, we develop a new method termed the PHENOTYPE method. In this method, the prior probabilities of phenotypes and the transitional probability from a phenotype to another phenotype will be used to establish various likelihood formulas.

## Phenotype method

We begin our discussion with the symbol  $\mathcal{G} \triangleright \mathcal{P}$ , whose meaning is that  $\mathcal{G}$  is a genotype determining the phenotype  $\mathcal{P}$ , i.e.  $\mathcal{G} \supseteq \mathcal{P}$  and  $\forall A \in \mathcal{G} \rightarrow A \in \mathcal{P}$ , where  $\supseteq$  is the inclusion of multisets. If the null alleles (e.g.  $A_j$ ) are considered, the conditions should be

revised to  $\mathcal{G} \supseteq \mathcal{P}$  and  $\forall A \in \mathcal{G} \rightarrow A \in \mathcal{P} \cup \{A_y\}$ . Under the revised conditions, our models will accommodate the effect of null alleles.

The formulas of transitional probabilities  $T(\mathcal{P}_O | \mathcal{P}_F)$  and  $T(\mathcal{P}_O | \mathcal{P}_F, \mathcal{P}_M)$  are first established, whose expressions are:

$$T(\mathcal{P}_O | \mathcal{P}_F) = \sum_{\mathcal{G}_F \in \mathcal{P}_F} \sum_{\mathcal{G}_O \in \mathcal{P}_O} \Pr(\mathcal{G}_F | \mathcal{P}_F) T(\mathcal{G}_O | \mathcal{G}_F) T(\mathcal{P}_O | \mathcal{G}_O), \quad (3)$$

$$T(\mathcal{P}_O | \mathcal{P}_F, \mathcal{P}_M) = \sum_{\mathcal{G}_F \in \mathcal{P}_F} \sum_{\mathcal{G}_M \in \mathcal{P}_M} \sum_{\mathcal{G}_O \in \mathcal{P}_O} \Pr(\mathcal{G}_F | \mathcal{P}_F) \Pr(\mathcal{G}_M | \mathcal{P}_M) T(\mathcal{G}_O | \mathcal{G}_F, \mathcal{G}_M) T(\mathcal{P}_O | \mathcal{G}_O), \quad (4)$$

where  $\mathcal{G}_F$  ( $\mathcal{G}_M$  or  $\mathcal{G}_O$ ) is taken from all genotypes determining  $\mathcal{P}_F$  ( $\mathcal{P}_M$  or  $\mathcal{P}_O$ );  $\Pr(\mathcal{G}_F | \mathcal{P}_F)$  and  $\Pr(\mathcal{G}_M | \mathcal{P}_M)$  are two posterior probabilities, which can be calculated by the Bayes formula:

$$\Pr(\mathcal{G} | \mathcal{P}) = \frac{T(\mathcal{P} | \mathcal{G}) \Pr(\mathcal{G})}{\Pr(\mathcal{P})};$$

and  $T(\mathcal{P}_O | \mathcal{G}_O)$  is the transitional probability from  $\mathcal{G}_O$  to  $\mathcal{P}_O$ , which is calculated by:

$$T(\mathcal{P} | \mathcal{G}) = \mathcal{B}_{\mathcal{P}=\emptyset} \beta + \mathcal{B}_{\mathcal{G} \in \mathcal{P}} (1 - \beta),$$

in which  $\beta$  is the negative amplification rate, and  $\mathcal{P} = \emptyset$  means that  $\mathcal{P}$  is a negative phenotype (it may be caused by either a null allele homozygote or a negative amplification).

Because each genotype may encounter an amplification failure, the candidate genotypes determining a negative phenotype at a locus are, strictly speaking, all possible genotypes at this locus. This will create a problem for the calculations of the transitional probabilities. This is because there are up to  $\binom{v+K-1}{v}$  genotypes at a locus, where  $v$  is the ploidy level and  $K$  is the number of alleles at this locus. For example, the number of genotypes at an octo-allelic locus for tetrasomic (hexasomic, octosomic or decasomic) inheritance is up to 330 (1716, 6435, or 19,448). For this reason, we do not consider the candidate genotypes determining any negative phenotypes. In other words, all negative phenotypes are discarded in the polysomic inheritance model during the analytical process. However, they will still be used in the allele frequency estimation so as to estimate the negative amplification rate  $\beta$  and the null allele frequency  $p_y$ .

Next, the likelihood formulas for all three categories are established. For example, if self-fertilization is not considered, the likelihoods  $\mathcal{L}(H_1)$  and  $\mathcal{L}(H_2)$  for the first category can be simply obtained by replacing  $\mathcal{G}_A$  with  $\mathcal{P}_A$  and  $\mathcal{G}_O$  with  $\mathcal{P}_O$  in Equation (2), whose expressions are:

$$\begin{aligned} \mathcal{L}(H_1) &= \Pr(\mathcal{P}_A) [(1-e)^2 T(\mathcal{P}_O | \mathcal{P}_A) + 2e(1-e) \Pr(\mathcal{P}_O) + e^2 \Pr(\mathcal{P}_O)], \\ \mathcal{L}(H_2) &= \Pr(\mathcal{P}_A) \Pr(\mathcal{P}_O), \end{aligned}$$

where  $\Pr(\mathcal{P}_A)$  and  $\Pr(\mathcal{P}_O)$  are, respectively, the frequencies of  $\mathcal{P}_A$  and  $\mathcal{P}_O$ , which can be calculated by Equation (A5), and the transitional probability  $T(\mathcal{P}_O | \mathcal{P}_A)$  is calculated by replacing  $\mathcal{P}_F$  with  $\mathcal{P}_A$  in Equation (3), i.e.  $T(\mathcal{P}_O | \mathcal{P}_A) = T(\mathcal{P}_O | \mathcal{P}_F = \mathcal{P}_A)$ . The likelihood formulas for each category under the condition of either self-fertilization or not are given in Supplementary Appendix E.

## Estimation of genotyping error rate

For a genotypic dataset, it is mathematically impossible to estimate the genotyping error rate  $e$  without any additional

information (e.g. the information of pedigree or replication). We will develop a genotyping error rate estimator based on the pedigree data, including the known parents and the identified parents (at a high confidence level, e.g. 99%). We refer to a parent-offspring pair extracted from the pedigree data as a *reference pair*, and a father-mother-offspring trio as a *reference trio*.

For genotypic data, we assume that the allelic dosage is known so there are no null alleles. For the phenotypic input, all candidate genotypes and their gametes will be extracted, including the genotypes with null alleles, and the pair (or trio) mismatch is identified by whether the parent (or the parents) is able to produce the offspring (see Supplementary Appendix I for details). Therefore, each mismatch in our models can only be caused by genotyping errors or the false parent(s). Pair mismatches can be used in all three categories, but trio mismatches can only be used in the second and the third categories. In this section, we will use pair mismatches to describe how to estimate the genotyping error rate.

Let  $\delta$  be the probability of observing a pair mismatch in a true parent-offspring pair under the condition that any individual has encountered a genotyping error. In our genotyping error model,  $\delta$  is equal to the exclusion rate for the first category, i.e. the probability that two random genotypes are mismatched. We do not estimate  $\delta$  by simulation or by allele frequencies because those approaches can be influenced by the errors in the estimated parameters. Instead, we directly estimate  $\delta$  from the input genotypes/phenotypes with a Monte-Carlo algorithm, whose procedures are broadly as follows: randomly sample a large number of individual pairs from the input samples with replacement, and then treat each as a parent-offspring pair, and finally calculate the probability that their genotypes/phenotypes at a locus are mismatched, which is used as  $\hat{\delta}$  at this locus.

Let  $\gamma$  be the probability of observing a pair mismatch in a true parent-offspring pair. Since each mismatch observed in the true parent-offspring pairs can only be caused by the genotyping error, if we denote  $E$  for  $1 - (1 - e)^2$ , then  $\gamma = E\delta$ . Noticing that the estimate  $\hat{\gamma}$  can be calculated from the reference pairs in a single application or in all available applications based on the same dataset, the single-locus estimate  $\hat{E}_1$  of  $E$  at the  $1^{\text{th}}$  locus can be expressed as  $\hat{E}_1 = \hat{\gamma}_1 / \hat{\delta}_1$ .

If we assume that there are  $n_{r1}$  reference pairs at the  $1^{\text{th}}$  locus and that  $n_{m1}$  is the number of pair mismatches in these reference pairs, then  $n_{m1}$  as a random variable obeys the binomial distribution  $B(n_{r1}, \gamma)$ , so  $\text{Var}(n_{m1}) = n_{r1}\gamma(1 - \gamma)$ . Because  $1 - \gamma$  is close to one, the variance  $\text{Var}(\hat{\gamma}_1)$  can be approximately expressed as  $\text{Var}(\hat{\gamma}_1) \approx \gamma_1/n_{r1}$ . Because  $\hat{E}_1 = \hat{\gamma}_1/\hat{\delta}_1$  and  $\gamma = E\delta$ , then  $\text{Var}(\hat{E}_1 \hat{\delta}_1) \approx (E\delta_1)/n_{r1}$ . Now, by substituting  $\delta_1$  with  $\hat{\delta}_1$ , it follows that  $\text{Var}(\hat{E}_1) \approx E/(n_{r1}\hat{\delta}_1)$ . To minimize the variance of  $\text{Var}(\hat{E})$ , the inverse of  $\text{Var}(\hat{E}_1)$  can be used as the weight to calculate the multi-locus estimate  $\hat{E}$ . The unified weight  $w_1$  is therefore equal to  $n_{r1}\hat{\delta}_1/(\sum_{i=1}^p n_{r1}\hat{\delta}_i)$ , and  $\hat{E} = \sum_i w_i \hat{E}_i$ . Because the loci are unlinked, we have  $\text{Var}(\hat{E}) = \sum_i w_i^2 \text{Var}(\hat{E}_i)$ , hence  $\text{Var}(\hat{E}) \approx E/(\sum_i n_{r1}\hat{\delta}_i)$ .

The genotyping error rate  $e$  can now be estimated by the formula  $\hat{e} = 1 - \sqrt{1 - \hat{E}}$ . Moreover, because  $e \approx E/2$ , the variance  $\text{Var}(\hat{e})$  can be approximately expressed as  $\text{Var}(\hat{e}) \approx e/(2 \sum_i n_{r1}\hat{\delta}_i)$ . As described above, the inverse of  $\text{Var}(\hat{e})$  can be used to weight  $\hat{e}$  in multiple applications and datasets.

When the polyploid phenotypes are used, pair mismatches will be rare. Specifically, they are rare for the first category, because the single-locus exclusion rate is low (e.g. 0.01 for the hexaploid phenotypes at a hexa-allelic locus). Therefore, it is inaccurate to estimate  $e$  by pair mismatches. Relative to the first category, the single-locus exclusion rate for the second or the

third categories is high (e.g. 0.27 for the hexaploid phenotypes at a hexa-allelic locus). Hence, we can use trio mismatches to reliably estimate the genotyping error rates for the second and the third categories, and the details are described in Supplementary Appendix F.

### Estimation of sample rate

For an individual offspring, the probability that one of its true parents is sampled is defined as the *sample rate*, denoted by  $p_s$ . The probability that an alleged parent (or a pair of alleged parents) of an offspring is assigned at a confidence level is called the *assignment rate*, denoted by  $a$ . Specifically, we denote  $a_c$  for the assignment rate when the true parent(s) is sampled, and  $a_u$  for the assignment rate when the true parent(s) is not sampled. Therefore,  $a$  is a weighted average of  $a_c$  and  $a_u$ .

We now develop a simple but robust estimator to estimate the sample rate from the assignment rate and begin our discussion with how to estimate the sample rate by using one application. For convenience, we will replace “the father” with “one parent” and “the mother” with “the other parent” in the first and the second categories in a parentage analysis.

For the first and the second categories, we have  $a = p_s a_c + (1 - p_s) a_u$ , so  $p_s$  can be estimated by:

$$\hat{p}_s = \frac{\hat{a} - \hat{a}_u}{\hat{a}_c - \hat{a}_u}. \quad (5)$$

For the third category, if the sexes are known, then  $a = p_s^2 a_c + (1 - p_s^2) a_u$ , so  $p_s$  can be estimated by:

$$\hat{p}_s = \sqrt{\frac{\hat{a} - \hat{a}_u}{\hat{a}_c - \hat{a}_u}}. \quad (6)$$

If the sexes are unknown, then  $a = p_c a_c + (1 - p_c) a_u$ , where  $p_c$  is the probability that the true parents are sampled, which can be expressed as  $p_c = s_u p_s + (1 - s_u) p_s^2$ , in which  $s_u$  is the proportion of selfed offspring in this application. Hence  $\hat{p}_c = \frac{\hat{a} - \hat{a}_u}{\hat{a}_c - \hat{a}_u}$ , and the sample rate  $p_s$  can be estimated by:

$$\hat{p}_s = \frac{\hat{s}_u - \sqrt{\hat{s}_u^2 + 4\hat{p}_c - 4\hat{s}_u\hat{p}_c}}{2\hat{s}_u - 2}. \quad (7)$$

The value of  $\hat{p}_s$  may be less than zero or greater than one. If this happens, we will truncate the value into the acceptable range [0, 1]. We will also set multiple confidence levels to estimate the selfing rate  $s_u$  for increased accuracy. For the situations of multiple applications and multiple confidence levels, the estimation of the sampling rate is shown in Supplementary Appendix G, along with the estimation of  $s_u$ .

### Evaluation

In this study, we use a computer simulation to create the genotypic and phenotypic datasets with disomic, tetrasomic or hexasomic inheritance, and then perform our parentage analysis by using these datasets. The performances of four methods under the same conditions are compared by four typical applications, where one method is the PHENOTYPE method, and the others are named the DOMINANT method (Rodzen et al. 2004) (named after the pseudo-dominant data used in this method), the SIBSHIP method (Wang 2016) (originating from the application “sibship reconstruction”) and the EXCLUSION method (Zwart et al. 2016). The accuracies of these four methods under natural conditions are

tested with an empirical microsatellite dataset for the highbush blueberry (Huber 2016). In addition, the performances of the genotyping error rate estimation and the sample rate estimation are also evaluated using the simulated datasets.

Both the DOMINANT and the SIBSHIP methods rely on first transforming the polyploid codominant phenotypic data into pseudo-dominant data. The same procedure as Kalinowski et al. (2007) is used for the DOMINANT method, and the likelihood formulas under this method are listed in Supplementary Appendix H, whose derivations are given by Gerber et al. (2000). Under the SIBSHIP method, a simulated-annealing algorithm is used to find the classification of optimal full-sib (or half-sib) families for the whole dataset by maximizing the likelihood, which is implemented in the software package COLONY (Wang and Scribner 2014). Under the EXCLUSION method, the effects of double-reduction and null alleles are incorporated, and the details of this method are described in Supplementary Appendix I.

### Simulated data

In order to evaluate these methods, we create some theoretical monoecious populations, each consisting only of individuals with disomic to decasomic inheritance for the genotypic data or disomic to hexasomic inheritance for the phenotypic data. We assumed that the population under scrutiny is genotyped at  $L$  unlinked loci under the PES model (Huang et al. 2019). The number of loci  $L$  is set from three to 12 (genotypes) or three to 18 (phenotypes) at an interval of three. The distance (in centimorgans) between each of these loci and its corresponding centromere is drawn from the uniform distribution  $U(0, 100)$ . The single chromatid recombination rate  $r_s$  is obtained by Haldane’s mapping function. Each locus is located with six amplifiable alleles that have uniform initial frequencies, with the initial null allele frequency set as 0.1 for the phenotypic data. For the genotypic data, null alleles are not simulated because the dosage of alleles within each genotype is known.

Huang et al. (2019) derived the genotypic frequencies under each of the four double-reduction models listed in Supplementary Table S1. However, the analytical solution of genotypic frequencies under inbreeding/selfing and double-reduction is still unknown. As an alternative, we give an approximated solution in Supplementary Appendix C by using the inbreeding coefficient  $F$  as an intermediate variable with the assumption that any inbreeding is only caused by self-fertilization. With this approximation, we generate the genotypes of the founder generation by Equation (A4). In order to let the genotypic frequencies reach their equilibrium state and avoid severe genetic drift, 2000 individuals are generated for the founder generation, and the population is allowed to reproduce for ten generations, each generation consisting of 2000 individuals.

During reproduction, the parents of each offspring are either two distinct individuals randomly chosen from the previous generation at a probability of  $1 - s$ , or the same individual (for self-fertilization) randomly chosen from the previous generation at a probability of  $s$ . The selfing rate  $s$  is set as three levels (0, 0.1 and 0.3). The following three procedures are designed to simulate meiosis: (1) the chromosomes are randomly paired and the alleles are exchanged between the pairing chromosomes at a probability of  $r_s$ ; (2) the chromosomes are randomly segregated into two secondary oocytes; and (3) the alleles within a chromosome are randomly segregated into two gametes. Fertilization is then simulated by the merging of two gametes.

Next, we reproduce two additional generations, each consisting of 100 individuals, to be used as the parents and offspring for

the subsequent analyses. To simulate the missing parents, 90% of parents and all offspring are sampled. To simulate the genotyping errors, each genotype is swapped with the genotype of another individual at the same locus at a probability of  $\frac{1}{2}e$  (where  $e$  is set as 0.01). To simulate negative amplification, each genotype is randomly set as  $\emptyset$  at a probability of  $\beta$  (where  $\beta$  is set as 0.05). The phenotypes are obtained by removing both the null and the duplicated alleles within genotypes. Then the generated genotypic (or phenotypic) dataset is used to perform the parentage analysis. The allele frequency estimation is described in Supplementary Appendix J.

For the first two categories in a parentage analysis, each is designated its own application [named Application (i) or (ii)]. Application (iii) refers to a third category in which the alleged fathers and the alleged mothers are drawn from two different collections (representing that the sexes are known). Application (iv) also refers to the third category in which the alleged fathers and the alleged mothers are drawn from the same collection (representing that the sexes are unknown).

In Application (i), for each of the 100 offspring, 89 individuals from the parental generation are used as alleged fathers. Application (ii) is performed for the offspring with their mother sampled. In this application, for each offspring, the true mother is known, and 89 individuals from the parental generation are used as the alleged fathers. For Applications (i) and (ii), the alleged fathers will include the true father if sampled but will exclude the true mother (except when the offspring is the product of self-fertilisation) to avoid interference. In Application (iii), for each offspring, 45 individuals (including the true father if sampled) from the parental generation are considered as the alleged fathers, with the remaining 45 individuals (including the true mother if sampled) as the alleged mothers. In Application (iv), for each offspring, all 90 individuals in the parental generation are considered as the alleged parents. We perform 100 replications for each configuration and calculate the average correct assignment rate for each configuration. Here, a configuration is defined as a combination of the parameters  $v$ ,  $L$  and  $s$ . A *correct assignment* therefore means that the true parents have been assigned and the value of  $\Delta$  is higher than the corresponding threshold.

For the PHENOTYPE method, there are many models to estimate the allele frequencies and the related parameters, and the ideal way is to try each and then choose the optimal one with the smallest Bayesian information criterion (BIC) (as in Huang et al. 2020). However, it is time consuming to evaluate each of them in each simulation. As an alternative, we choose two models that work well in most situations:  $PES_{0.25} + p_y + \beta + s$  for the phenotypic data and  $PES_{0.25} + \beta + s$  for the genotypic data. They denote the PES models with  $r_s = 0.25$  together with the considerations of null alleles (for phenotypes only), negative amplification and self-fertilization. Because the estimations of genotyping error rate  $e$  and sample rate  $p_s$  depend on the number of assigned parents, the performance of a less efficient method will be reduced again due to the inaccurate estimations of  $e$  and  $p_s$ . As the aim of our simulation is to evaluate the performance of four methods, not the influence of the estimations of  $e$  and  $p_s$ , the true values of  $e$  and  $p_s$  are used as the *a priori* information. We perform 2000 Monte-Carlo simulations to obtain various critical values of the statistic  $\Delta$ , and the correct assignment rates under three critical values ( $0$ ,  $\Delta_{0.8}$  and  $\Delta_{0.95}$ ) are recorded.

For both the DOMINANT and the SIBSHIP methods, the frequency  $p_{\text{dom}}$  of the dominant allele at a pseudo-dominant marker is estimated by  $\hat{p}_{\text{dom}} = 1 - \sqrt{1 - \hat{p}_{\text{tar}}}$ , where  $\hat{p}_{\text{tar}}$  is the observed probability that a randomly sampled phenotype contains the target

allele. For the DOMINANT method, we implement the calculations of likelihood formulas listed in Supplementary Appendix H in our simulation program. We also perform 2000 Monte-Carlo simulations to obtain the thresholds of  $\Delta$ , and record the correct assignment rates under the same thresholds as above. For the SIBSHIP method, we write the pseudo-dominant phenotypes, the allele frequency estimates and other necessary parameters into a COLONY V2.0.6.5 input file. To avoid interference by the other cases, a unique input file for each case is generated. After calling colony2p.exe by a command-line mode, the results can be read from the output files. The probability of the identified parent(s) is used as a confidence level to compare with the PHENOTYPE and DOMINANT methods. The EXCLUSION method is implemented in our simulation program. In this method, the alleged parent (or parent-pair) with the fewest mismatches is assigned. If multiple alleged parents (or parent-pairs) have the same number of mismatches, none of them is assigned. For this method, any confidence level is unavailable.

### Empirical data

We used a microsatellite dataset from the highbush blueberry (*Vaccinium corymbosum*) (Chapter 5, Huber 2016) to test the same four methods. The highbush blueberry has tetrasomic inheritance with no evidence of fixed heterozygosity (that indicates disomic inheritance; Krebs and Hancock 1989).

The blueberry samples were collected from Agriculture Agri-Food Canada blueberry plots in Abbotsford and Agassiz, BC., Canada (Huber 2016). Five controlled crosses, each with 25 to 30 offspring, were collected, resulting in a collection of 150 individuals, 143 of which were offspring. All samples were successfully amplified at 15 microsatellite loci, with the number of alleles sampled ranging from three to ten (mean  $\pm$  SD is  $5.60 \pm 2.33$ ).

There are altogether seven parents in these five controlled crosses. To increase the difficulty of our analysis, we also add 120 simulated false parents which are generated by randomly copying the phenotypes from the real individuals. Following the four applications for the simulated data, we designed four similar applications for these empirical data. Application (I) or (II) refers to identifying the father when the mother is either unknown or known. There are 286 cases (twice the number of offspring) for each application, and each case has either 60 alleged fathers (including one true parent and 59 false parents) for Application (I) or the known mother together with 60 alleged fathers (including one true parent and 59 false parents) for Application (II). Application (III) refers to identifying the father and the mother jointly in which the alleged fathers and the alleged mothers are drawn from two different collections. There are 143 cases for this application, each of which has 30 alleged fathers (including the true father and 29 false fathers) and 30 alleged mothers (including the true mother and 29 false mothers). Application (IV) refers to identifying the father and the mother jointly in which the alleged fathers and the alleged mothers are drawn from the same collection. There are also 143 cases for this application, each of which has 60 alleged parents of unknown sex (including two true parents and 58 false parents). The false parents used in each case are randomly sampled from a total of 125 false parents (including 120 simulated individuals and 5 natural individuals).

We randomly sample five to 15 loci from the dataset. For each value of  $L$ , 100 datasets are generated, each including 150 true individuals and 120 false parents. These datasets will be used to perform our parentage analysis by using the same four methods as described in the previous section. The analytical procedures are also the same as in the previous section except that the

number of Monte-Carlo simulations to obtain the thresholds of  $\Delta$  is 10,000 instead of 2000. The correct assignment rate will be used to measure the accuracy of each model.

### Evaluation of genotyping error rate and sample rate

We use the simulated data to evaluate the performances of both estimators for the genotyping error rate and the sample rate. The same four applications are used as previously described, and are still referred to as Applications (i) to (iv). We estimate the genotyping error rate and the sample rate for each application. Two pairs of sampling and genotyping conditions, *poor* and *good*, are selected, which are  $e=0.1$  and  $p_s=0.5$  for *poor*, or  $e=0.02$  and  $p_s=0.8$  for *good*. The remaining parameters are almost the same as those in the section *Simulated data*, in which  $s=0.1$ ,  $p_y=0.1$  and  $L$  is taken from 6 to 24 at an interval of three. We then perform 100 simulations for each configuration. The *PHENOTYPE* method is used to perform the parentage analysis with *a priori* genotyping error rate  $e=0.01$  and sample rate  $p_s=0.9$ . The allele frequencies are estimated under the  $PES_{0.25} + p_y + \beta + s$  model. The performances of both estimators are evaluated by the RMSE.

For the estimation of the genotyping error rate, the identified pairs (or trios) with a confidence level of 99% are considered as the reference pairs (or trios), with  $\delta$  estimated by randomly sampling 10,000 pairs (or trios). In Application (i),  $\hat{e}$  is estimated from the pair mismatch, whilst for the remaining applications  $\hat{e}$  is estimated from the pair or the trio mismatches.

For the estimation of sample rate, we use the weighted average of  $\hat{p}_s$  across three confidence levels (80%, 95% and 99%) for each application. Because  $\hat{a}_c$  and  $\hat{a}_u$  are obtained from the simulation, they may be influenced by any inaccurate simulation parameters, such as the sample rate, the selfing rate and the genotyping error rate. To improve the accuracy of these simulation parameters, we perform two rounds of analyses. The estimated sample rate and genotyping error rate in the first round are used as the *a priori* values in the second round. The results of the second round are used to evaluate the performance.

### Data availability

Supplementary material S1: The appendices and supplementary figures.

Supplementary material S2: The simulation parameters, output files, description of I/O format, figure plotting script and empirical dataset.

POLYGENE is written in C++ and C#, whose executables (Windows, Ubuntu and Mac OS X), source code and user manual are available on GitHub (<http://github.com/huangkang1987/polygene>). The simulation functions are “private void SIM\_PARENT1 ()” to “private void SIM\_PARENT3 ()” in “Form1.cs.”

Supplementary material is available at <https://doi.org/10.25387/g3.13272623>.

## Results

### Simulated data

For the four applications, each correct assignment rate as a function of  $L$  is denoted by a section of the overlapped bar charts, shown in [Figure 1](#) for the genotypic data or in [Figure 2](#), and Supplementary Figures S1 and S2 for the phenotypic data.

For the genotypic data, it can be seen from [Figure 1](#) that each correct assignment rate increases as the number of loci  $L$  also increases, whose values reach a steady state if  $L$  is large enough [e.g.  $L \geq 12$  for Application (i) or  $L \geq 9$  for the other applications]. The correct assignment rate generally reduces as the ploidy level

increases. Moreover, as the selfing rate increases the correct assignment rate also increases but the difference among different ploidy levels decreases.

For the phenotypic data, it can be seen from [Figure 2](#) that the correct assignment rate reduces as the ploidy level increases. The *PHENOTYPE* method outperforms the other methods, whose correct assignment rate at  $L=9$  is roughly the same as those of the other methods at  $L=18$ , indicating that the *PHENOTYPE* method can reduce the number of loci needed to achieve the same accuracy by 40% to 60%. This method is also less sensitive to changes in the ploidy level, but an additional 23% and 45% loci are still required to reach the same correct assignment rate in tetraploids and hexaploids, respectively.

Compared with the *DOMINANT* method, the performance of the *SIBSHIP* method is improved in Applications (i) and (ii) at a high  $L$  ( $\geq 15$ ), but is inferior in the other scenarios. The performance of the *EXCLUSION* methods is good in Applications (ii) to (iv) at a high  $L$  ( $\geq 15$ ) but is inapplicable in Application (i).

It can be seen from Supplementary Figures S2 and S3 that, like the results of genotypic data, the correct assignment rate is increased under most situations if the selfing rate is increased from 0 to 0.3. The assignment rate is reduced in Applications (ii) to (iv) under both the *SIBSHIP* and the *EXCLUSION* methods.

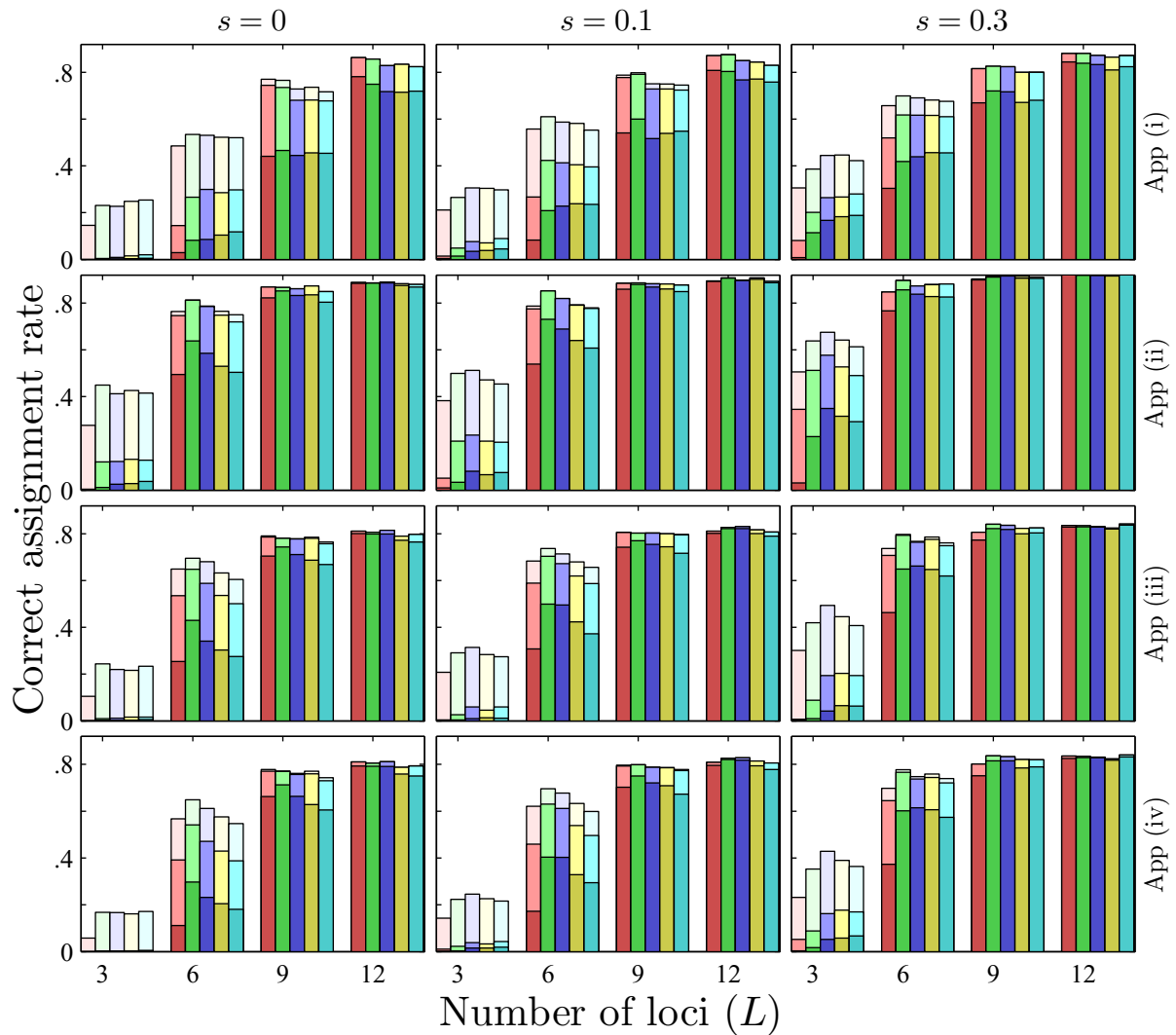
### Empirical data

The parentage assignment results from using each of the four methods and applying the phenotypic dataset of [Huber \(2016\)](#) are shown in [Figure 3](#). The results patterns are similar to those obtained from the simulated data. The *PHENOTYPE* method still outperforms the other three methods but to a lesser degree than when the simulated dataset was used, but the *PHENOTYPE* method can still achieve the same accuracy with only 75% of the loci needed for the other methods. The *EXCLUSION* method is still inaccurate and cannot be applied to real data in Application (I), but its performance is relatively good for the other applications when  $L > 10$ . The *DOMINANT* method performs worse than the other three methods for Application (IV), as does the *SIBSHIP* method for Application (I).

### Evaluation of genotyping error rate and sample rate

The results under both poor and good conditions are shown in [Figure 4](#) and Supplementary Figure S4, respectively. For the estimation of the genotyping error rate, it can be seen that the results are good due to the RMSE being reduced to a low level. For example, the RMSE at  $L=24$  is able to reach 0.02 in poor conditions or 0.005 in good conditions. The RMSE for Application (i) performs worse than for the other applications, and increases greatly as the ploidy level also increases. This is because only the pair mismatch can be used for this application, and the single-locus exclusion rate for the first category is small. The RMSE for Application (ii) performs better than for the other applications, because both the pair and the trio mismatches are used for this application, and the single-locus exclusion rate for the second category is usually higher than the other applications. The RMSE curves of Applications (iii) and (iv) are similar.

For the estimation of the sample rate, [Figure 4](#) and Supplementary Figure S4 show that the results are inferior to those for the estimation of the genotyping error rate. For example, the RMSE at  $L=24$  is only able to reach 0.05 in poor conditions or 0.02 in good conditions. Unlike the estimation of the genotyping error rate, the results for Application (i) are not obviously inferior to those for the other applications. This is because the assignment rate rather than the reference pairs is used to



**Figure 1** Correct assignment rate as a function of the number of loci  $L$  by using the genotypic data. Each row is designated an application and each column shows the simulation results for a different rate of selfing. Every correct assignment rate is denoted by a section of overlapping bar charts. The results of disomic to decasomic inheritances are shown by red, green, blue, yellow and azure bars, respectively. The bars with light, medium and bright colors denote in turn the correct assignment rates with the thresholds  $0$ ,  $\Delta_{0.80}$  and  $\Delta_{0.95}$ . Applications (i) to (iv) correspond to (i) identifying one parent when the other is unknown, (ii) identifying one parent when the other is known, (iii) identifying parents of known sexes jointly, and (iv) identifying parents jointly with unknown sexes.

estimate the sample rate, causing the results influenced less by the low single-locus exclusion rate.

The results for Application (ii) are poorer than those for estimating the genotyping error rate because fewer cases ( $\approx 50$  cases) are used (about half of the true mothers are not sampled). If Applications (i) and (ii) use the same number of cases, then the performance of Application (ii) would be better than Application (i). Because Application (ii) also uses the mother's data, which can better distinguish the true and the false fathers, the difference between  $a_c$  and  $a_u$  in Application (ii) is larger than that in Application (i) under the same conditions (*e.g.* Supplementary Figure S3).

The results for Application (iii) are usually better than those for the other applications. This is because Application (iii) does not need to estimate the selfing rate and has a larger sample size (100 cases). However, the selfing rate has to be estimated for Application (iv), and thus the results are less accurate than for Application (iii).

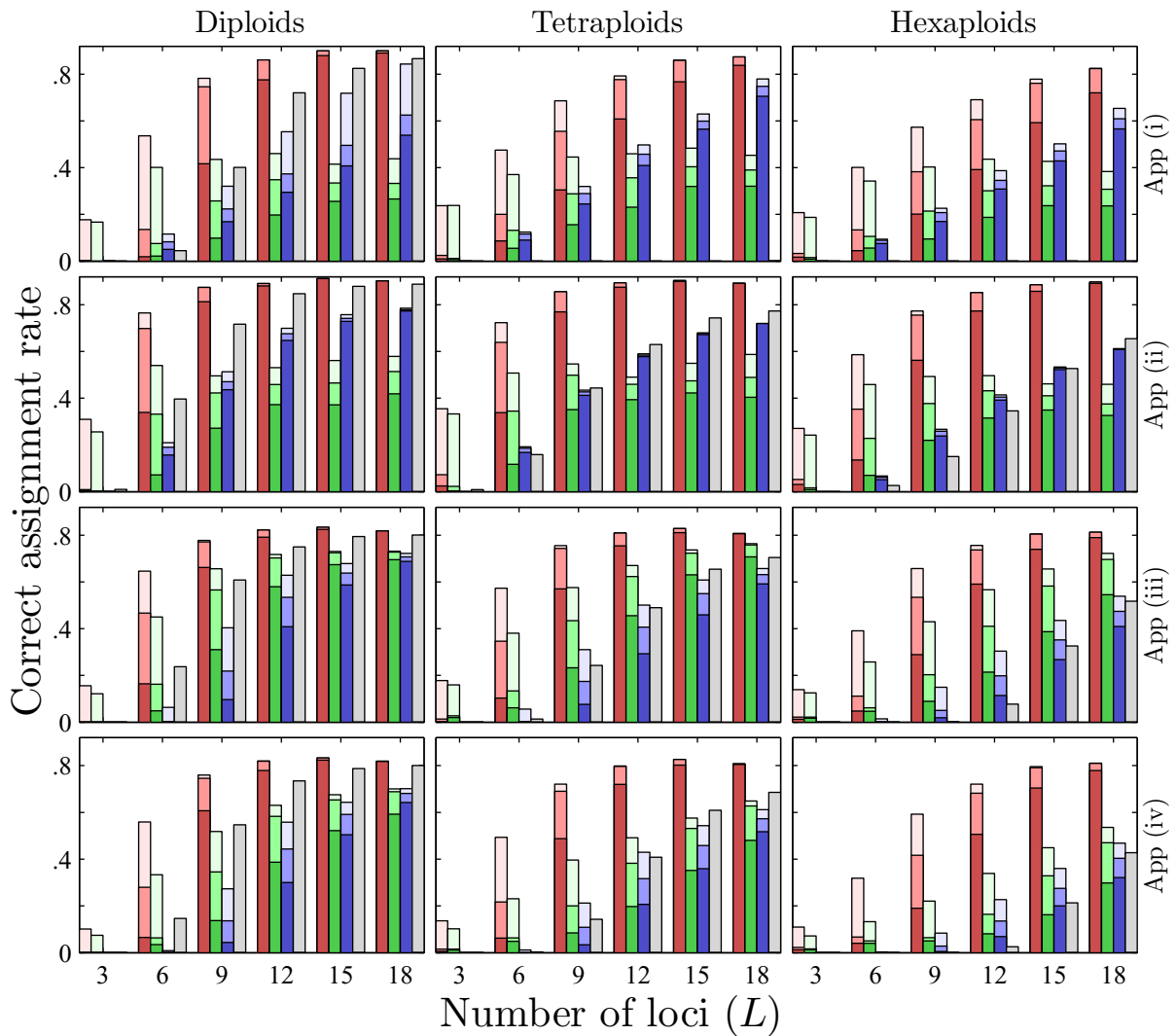
## Discussion

### Inheritance model

Meiosis in polyploids is complex. Disomic and polysomic inheritances are two extremes, and many autopolyploid taxa represent the intermediate stages (Butruille and Boiteux 2000). Allopolyploids (such as the segmental allopolyploids) can also display intermediate inheritance at some loci (Stift *et al.* 2008). In addition, some autopolyploid species can also form bivalent, univalent and other types of valents during meiosis (Lloyd and Bomblies 2016). The formation of different types of valents may influence the sterility of the gametes or the seeds (Solís Neffa and Fernández 2000).

For the autopolyploids with pure disomic inheritance, we can adopt the RCS model to simulate disomic inheritance. This is because the genotypic frequencies, gamete frequencies and transitional probabilities in the RCS model are the same as those for disomic inheritance. These probabilities are of interest for parentage analysis. The difference between the RCS model and





**Figure 2** Correct assignment rates as a function of the number of loci  $L$  by using the phenotypic data at a selfing rate of 0.1. Each row is designated an application and each column shows the simulation results for a different ploidy level. The results for the `PHENOTYPE`, `DOMINANT`, `SIBSHIP` and `EXCLUSION` methods are shown by the red, green, blue and gray bars, respectively. The bars with light, medium and bright colors denote in turn the correct assignment rates with the confidence levels 0, 80% and 95%.

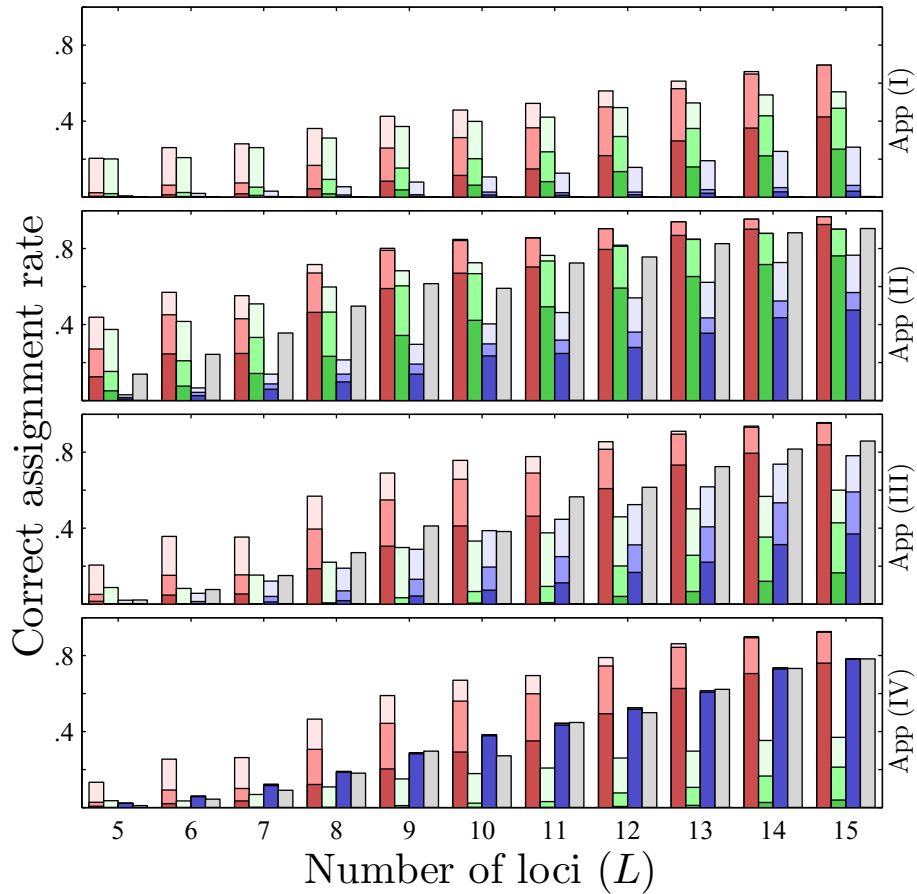
disomic inheritance is that 100% multivalent formation is assumed in the former, whilst 100% bivalent formation is assumed in the latter. For the allopolyploids with pure disomic inheritance, all diploid methods including those of parentage analysis can be used if the genotypes at different isoloci are identified.

For intermediate inheritance, *e.g.* 50% bivalent and 50% multivalent gamete formation, regardless of how complex the nature of meiosis, identical-by-double-reduction (IBDR) alleles will be present in the resulting fertile gametes (Huang et al. 2019). For this reason, a generalized model was proposed, which uses  $\lfloor v/4 \rfloor$  double-reduction rates in the calculation of genotypic frequencies and is able to describe meiosis patterns including that for intermediate inheritance (Huang et al. 2019). However, this model is too complex because it has  $\lfloor v/4 \rfloor$  more degrees of freedom than the RCS (PRCS or CES) model. It is difficult to accurately estimate each double-reduction rate and thus is unrealistic to apply to many actual conditions. Even if these double-reduction rates are estimated, this model will often be suboptimal to other models because of the requirement for more degrees-of-freedom to explain various trends in a data set resulting in a higher BIC.

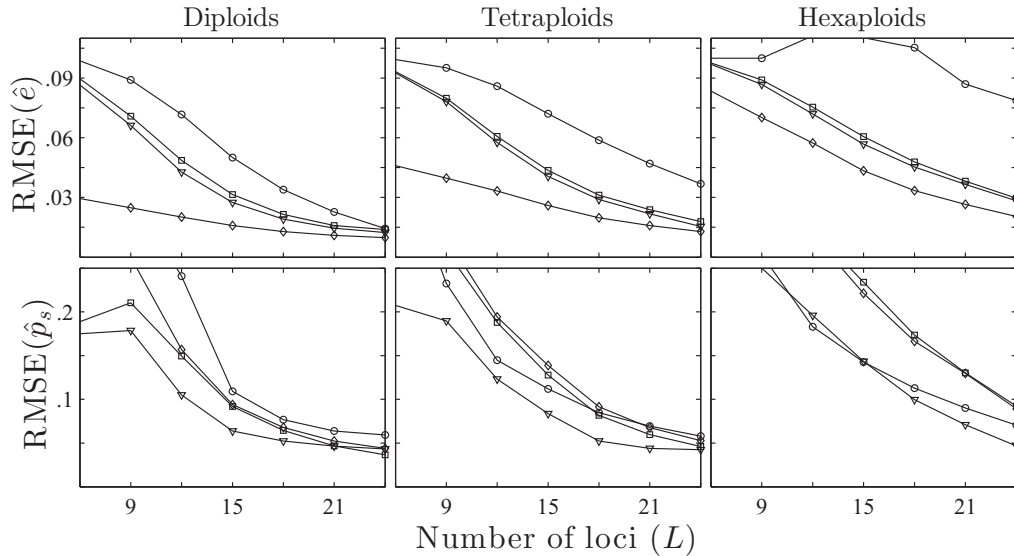
To better approximate the natural patterns, a simplified version of the generalized model was developed, named the PES model, which accommodates the single chromatid recombination rate  $r_s$  as an additional parameter to calculate the genotypic frequencies (Huang et al. 2019). Especially, this model is equivalent to either the RCS model if  $r_s = 0$ , or the CES model if  $r_s = 1$ . Our software provides three PES-related models, which are the PES0.25, the PES0.5 and the PES estimate  $r_s$ . The former two models do not increase their degrees-of-freedom because they use a fixed value of  $r_s$ . We suggest to evaluate candidate models by the BIC and chose the optimal model with the lowest BIC (as in Huang et al. 2020).

### Performance of parentage analysis

For the genotypic data, the results for polyploids are generally similar to those for diploids (Figure 1). The correct assignment rate tends to increase if the ploidy level ranges from two to four, whilst the assignment rate decreases with a ploidy level that ranges from four to ten. However, this trend is weakened as the selfing rate increases.



**Figure 3** Correct assignment rates as a function of the number of loci  $L$  by using the phenotypic dataset of Huber (2016). Each row denotes an application. The methods, confidence levels and the definitions of bars together with their shading are as for Figure 2.



**Figure 4** The RMSE of the estimated genotyping error rate  $\hat{e}$  or the estimated sample rate  $\hat{p}_s$  as a function of the number of loci  $L$  at  $e = 0.1$  and  $p_s = 0.5$ . Each column shows the results for a different ploidy level. The curves with circular, rhombic, triangular and squared markers denote the results for Applications (i), (ii), (iii), and (iv), respectively.

These phenomena have at least three not necessarily mutually exclusive explanations. (i) At a high polyploid level, a genotype has more allele copies and so contains more genetic information (Huang et al. 2014). This can improve the performance of parentage analysis and many other population genetics

analyses (e.g. the estimation of allele frequencies, genetic diversity,  $F$ -statistics, and relatedness coefficients). (ii) At a high polyploid level, the false parents are more likely to share the same alleles with the offspring, which may reduce the correct assignment rate. For example, if the ploidy level is high, reaching 1000,

the false parents will share the same alleles with the offspring at a hexa-allelic locus. This is similar to when biallelic loci are used in tetraploids or hexaploids, the details of which are discussed in the following section. (iii) Selfing is able to reduce the difference among ploidy levels and improve the performance of our parentage analysis. Each of these three explanations will also be reflected in the phenotype results and are described at the end of this section.

For the phenotypic data, the results for polyploids are generally inferior to those for diploids for each application and for each method (e.g. see Figure 2). The PHENOTYPE method performs best among all four methods, saving at least 25% more loci than the other methods (e.g. see Figures 2 and 3), whose performances are stable for all applications.

For the four applications, the results of the PHENOTYPE method for diploids (Figure 2, and Supplementary Figures S2 and S3) are slightly inferior to those for the genotypic data (Figure 1). This is because null alleles are simulated for the phenotypic data. In the absence of null alleles, each phenotype is only determined by one genotype for diploids. Therefore, both results under such condition are identical (data not shown).

For the DOMINANT (Rodzen et al. 2004) and SIBSHIP (Wang and Scribner 2014) methods, the results are suboptimal to those of the PHENOTYPE method (e.g. see Figure 3 and Supplementary Figure S2). In both the dominant and sibship methods, the polyploid co-dominant phenotypic data are transformed into the pseudo-dominant data, and the diploid procedures for a parentage analysis are subsequently used to perform an analysis. During transformation, genetic information is lost (Wang and Scribner 2014) and some noise is also introduced. For example, in the pseudo-dominant approach the pseudo-dominant loci are assumed to be unlinked. In fact, because there are at most  $v$  alleles in a phenotype, the presence of an allele in a phenotype will reduce the probability of observing the other alleles in this phenotype, and so these loci are negatively correlated rather than unlinked. In addition, for the pseudo-dominant approach, many factors that affect the parentage analysis are not considered, such as double-reduction, null alleles, negative amplification, and inbreeding/selfing.

The EXCLUSION method (Zwart et al. 2016) performs well in Applications (ii) to (iv), and the results are better than those for both the DOMINANT and the SIBSHIP methods but only if  $L$  is high (e.g. see Figure 3 and Supplementary Figure S3). However, the EXCLUSION method cannot be used for Application (i) because the single-locus exclusion rate in the first category is too low (e.g. 0.01 for hexaploid phenotypes at a hexa-allelic locus). Therefore, hundreds of loci are needed in order to exclude the false parents. This feature also influences the estimation of the genotyping error rate, such that the RMSE for Application (i) is highest (Figure 4).

From our simulation results, self-fertilization improves the accuracy of a parentage analysis, and reduces the variation of accuracies among different ploidy levels (Figures 1 and 2, and Supplementary Figures S2 and S3). This is because the genotypes become more homozygous as the selfing rate increases. If the selfing rate is one, all genotypes will become homozygotes at an equilibrium state. In such a case, each individual can be regarded as a haploid, and the ploidy level will not affect the accuracy of a parentage analysis.

## Genotyping error model

We follow Marshall et al. (1998) and consider the genotyping error as the replacement of a true genotype with a random genotype

according to the genotypic frequencies. This error model assumes that the genotyping error will not change the genotypic frequencies and therefore facilitates subsequent model development. However, actual genotyping errors can be complex, such as allelic dropouts, false alleles, mutations, miscalling, contaminated DNA and errors in data entry (Taberlet et al. 1996; Wang 2004). Although our error model may be unrealistic, it is likely the only applicable solution to our study for four reasons.

First, genotyping errors will interfere with the parentage analysis by reducing the LOD of the true parents. Assuming the expected LOD of an erroneously typed true parent is equal to the expected LOD of the false parents, then the genotyping error is equivalent to discarding  $eL$  locus. This estimation is conservative and overestimates the influence of genotyping errors. This is because some kinds of genotyping errors (e.g. allelic dropouts, false alleles) yield similar genotypes that share some identical alleles with the original genotypes. If there are sufficient loci (e.g. 18 microsatellites) and the genotyping error rate is not high (e.g.  $\epsilon < 0.05$ ), then the correct assignment rate still can be kept at a high level and the influence of genotyping errors are minimized for any error model. Second, the main function of considering genotyping errors is to tolerate some mismatches between the alleged parent and the offspring. Using an alternative error model may slightly change the LOD but will not change the parentage analysis results. For example, although Marshall et al. (1998) error model is incorrect (Kalinowski et al. 2007) the results of CERVUS v2.0 will still be applicable. Third, an alternative error model may not outperform the current error model. For any alternative model, additional parameters are required to describe the frequencies of different kinds of genotyping errors. These parameters can be assigned from *a priori* information, additional experiments (e.g. repeat genotyping) or additional estimators, and will increase the model complexity (i.e. degrees-of-freedom), Akaike information criterion (Akaike 1974) and Bayesian information criterion (Schwarz 1978). Finally, an alternative error model cannot be applied to our study due to difficulties in computation caused by increasing the number of alleles and reducing the calculation speed during simulation. Unfortunately, our current simulation speed is already slow even though we use several optimization methods (see the Optimization and complexity section).

## Genotyping error rate and sample rate

Our estimator for the genotyping error rate  $e$  is asymptotically unbiased as the number of loci increases. The bias of  $\hat{e}$  is from the estimation of  $\gamma$ . Because  $\gamma$  is estimated from any mismatches in the reference pairs or trios that are extracted from the identified parent(s), the confidence level of the true parents with few mismatches are successfully identified at a high probability. As a result, the value of  $\hat{\gamma}$  may be underestimated.

The estimation of the genotyping error rate does not use any simulation ( $\gamma$  is estimated from the reference pairs or trios, and  $\delta$  is estimated from the distribution of the observed genotypes/phenotypes). This means that the estimator is not only robust but also insensitive to any errors in the simulation parameters (such as the allele frequency, negative amplification rate, selfing rate, sample rate, or the genotyping error rate). Any errors in these simulation parameters can only slightly affect the identified parents, which will not significantly affect the accuracy of  $\hat{e}$ . However, this estimator needs sufficient loci to identify the reference pairs or trios. For instance, if  $e = 0.1$  and  $p_s = 0.5$ , at least 15 loci are required in order to estimate the genotyping error rate for hexaploids in Application (i) (Figure 4).

Compared with the genotyping error rate, the estimation of the sample rate  $p_s$  is less accurate and more sensitive to errors in the simulation parameters. There are at least three not necessarily mutually independent explanations for these patterns. (i) The estimate of the genotyping error rate is the weighted average of single-locus estimated values across all loci, where the actual sample size is  $\sum_i n_i$ . Whilst the sample rate is estimated only once for all loci, the actual sample size is the number of cases  $n_c$  (see Supplementary Appendix G). (ii) The sample rate estimator is biased in all categories in a parentage analysis because  $\hat{p}_s$  is truncated into the range  $[0, 1]$  and the operation of the square root is used in the third category. (iii) The simulation is used to obtain  $\hat{a}_c$  and  $\hat{a}_u$  for the estimation of the sample rate, whilst the parameters used for simulation may be inaccurate (e.g. a prior  $e$  and  $p_s$ ). Any errors in  $\hat{a}_c$  and  $\hat{a}_u$  can be passed to  $\hat{p}_s$ , but such errors can be eliminated by increasing the number of loci. When the number of loci are sufficient,  $\hat{a}_c$  will be close to one, and  $\hat{a}_u$  to zero. We suggest that users perform two rounds of estimation so as to reduce such errors as we have in the evaluation above.

### Polymorphism of loci

Because polyploids have more allele copies in a genotype, the false parents are more likely to share the same alleles with the offspring. Therefore, data resulting from the use of biallelic markers, e.g. single nucleotide polymorphism (SNPs), are unsuitable for performing a polyploid parentage analysis.

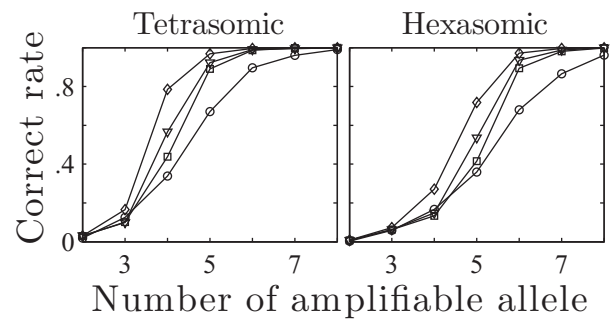
We will illustrate this by using the exclusion approach for the first category. For a given alleged parent, if its phenotype  $\mathcal{P}_A$  does not share any allele with its offspring phenotype  $\mathcal{P}_O$ , then it can be excluded as a true parent. If we assume that the double-reduction model is the RCS model, and that there are no interference factors (such as genotyping errors, self-fertilization, null alleles or negative amplification), then the exclusion rate  $\text{Excl}_1$  at a biallelic locus for the first category is:

$$\text{Excl}_1 = \Pr(\mathcal{P}_O = A, \mathcal{P}_A = B) + \Pr(\mathcal{P}_O = B, \mathcal{P}_A = A) = 0.5^{2v-1},$$

where  $A$  and  $B$  are the two alleles at this locus. The values of  $\text{Excl}_1$  from disomic to decasomic inheritances are in turn 0.125,  $7.813 \times 10^{-3}$ ,  $4.883 \times 10^{-4}$ ,  $3.052 \times 10^{-5}$  and  $1.907 \times 10^{-6}$ . This sequence decreases exponentially, indicating that the false parents become less likely to be excluded as the ploidy level increases. Moreover, the number of loci required to achieve the combined exclusion rate 0.95 is  $\ln(0.05)/\ln(1 - \text{Excl}_1)$ , whose values from disomic to decasomic inheritances are in turn 22, 382, 6134, 98,163, and 1,570,625.

Although next-generation sequencing (NGS) is able to segregate millions of SNPs, two reasons make it difficult to directly perform a parentage analysis with data obtained by using SNPs. First, the allele frequencies of most SNPs are not uniform, which reduces the exclusion rate. Second, adjacent SNPs are closely linked. This will reduce the accuracy of results because the genetic markers are assumed to be unlinked in all parentage analysis models.

Fortunately, haplotype assembly (Aguir and Istrail 2013), phased sequencing (Yang et al. 2011; Manching et al. 2017) and haplotype inference (Neigenfind et al. 2008) can all help to maintain the efficiency of NGS data, and can segregate multi-allelic loci by combining the closely linked variants so as to increase the single-locus polymorphism. Additionally, polyploid genotype calling can directly call back the genotypes but can currently



**Figure 5** Correct assignment rates as a function of the number of amplifiable alleles under the `PHENOTYPE` method. Twenty loci with uniform allele frequencies of amplifiable alleles are used. The threshold and the selfing rate are set as  $\Delta_{0.95}$  and 0.1, respectively. The remaining parameters and configurations are as for the simulated dataset. Each column shows the results for either tetrasomic or hexasomic inheritance. Each curve denotes the result for an application, whose definitions are as for Figure 4.

only be applied to the biallelic variants (Carley et al. 2017; Weiß et al. 2018).

Multi-allelic markers can also be influenced by the same problem. We perform a simple simulation to describe the influence of the number of amplifiable alleles on the correct assignment rate, in which 20 loci with uniform amplifiable allele frequencies are used to perform our parentage analysis under the `PHENOTYPE` method (Figure 5). The correct assignment rate is much increased if the number of amplifiable alleles equates broadly to the ploidy level  $v$ , indicating that to achieve the optimal result, the number of amplifiable alleles should be greater than or equal to  $v$  (Figure 5). More loci are required if loci with relatively low levels of polymorphism are used. We suggest therefore to use highly polymorphic loci to perform parentage analysis.

### Optimization and complexity

We use multi-threading, dynamic programming and genotype/phenotype indexing to optimize computational speed in the parentage analysis module in `POLYGENE`. The dynamic programming stores the likelihoods or LODs into a table so as to avoid repeated calculations. The genotype/phenotype indexing only records the hash values of genotypes/phenotypes for each individual, and the information of genotypes/phenotypes are saved in a hash table, that also includes the alleles, various frequencies (or prior/posterior probabilities), possible gametes and the number of occurrences.

All of these simulations took a total of three weeks to compute using a powerful workstation (Xeon E5 2699V4 36 cores). Computing efficiency will also be affected by the ploidy level  $v$  and the number of alleles  $K$  due to four main reasons: (i) the number of phenotypes increases as  $v$  and  $K$  increase, which reduces the efficiency of dynamic programming because more memory is required to store the likelihoods or LODs; (ii) the average number of genotypes determining a phenotype increases as  $v$  and  $K$  increase, which decelerates the calculation of likelihoods or LODs; (iii) the average number of gametes produced by a zygote increases as  $v$  and  $K$  increase, which decelerates the calculation of  $T(G_O|G_F)$  and  $T(G_O|G_F, G_M)$  in Equation (A6); (iv) the number of terms in Equation (A7) increases as  $v$  and  $K$  increase, which decelerates the calculation of  $T(g|G)$  in Equation (A7). These four factors collectively and multiplicatively increase the complexity of the calculations. It is therefore not possible to

perform an extensive simulation for highly polymorphic loci (e.g.  $K > 7$ ) or for high ploidy levels (e.g.  $v = 8$  or  $v = 10$ ).

## Acknowledgments

We would like to thank the subject editor and two anonymous reviewers for their suggestions and comments.

KR and BGL designed the project, KH and KR constructed the model, GH provided the data, KH wrote the draft, GH and DD edited the manuscript.

## Funding

This study is funded by the Strategic Priority Research Program of the Chinese Academy of Sciences (grant number XDB31020302), the National Natural Science Foundation of China (grant numbers 31730104, 31770411 and 31572278), the Innovation Capability Support Program of Shaanxi (2021KJXX-027), the Young Elite Scientists Sponsorship Program by CAST (grant number 2017QNRC001), the National Key Programme of Research and Development, Ministry of Science and Technology (grant number 2016YFC0503200), and the Shaanxi Science and Technology Innovation Team (grant number 2019TD-012). Derek W. Dunn is supported by Shaanxi Province Talents 100 Fellowship.

**Conflicts of interest:** The authors declare no conflict of interest.

## Literature cited

- Aguiar D, Istrail S. 2013. Haplotype assembly in polyploid genomes and identical by descent shared tracts. *Bioinformatics*. 29: i352–i360.
- Akaike H. 1974. A new look at the statistical model identification. *IEEE Trans Automat Contr*. 19:716–723.
- Barker MS, Arrigo N, Baniaga AE, Li Z, Levin DA. 2016. On the relative abundance of autopolyploids and allopolyploids. *New Phytol*. 210:391–398.
- Bezemer N, Krauss S, Phillips R, Roberts D, Hopper S. 2016. Paternity analysis reveals wide pollen dispersal and high multiple paternity in a small isolated population of the bird-pollinated *Eucalyptus caesia* (Myrtaceae). *Heredity*. 117:460–471.
- Blouin MS. 2003. DNA-based methods for pedigree reconstruction and kinship analysis in natural populations. *Trends Ecol Evol*. 18: 503–511.
- Butruille D, Boiteux L. 2000. Selection-mutation balance in polysomic tetraploids: impact of double reduction and gametophytic selection on the frequency and subchromosomal localization of deleterious mutations. *Proc Natl Acad Sci USA*. 97:6608–6613.
- Carley CAS, Coombs JJ, Douches DS, Bethke PC, Palta JP, et al. 2017. Automated tetraploid genotype calling by hierarchical clustering. *Theor Appl Genet*. 130:717–726.
- Darlington CD. 1929. Chromosome behaviour and structural hybridity in the Tradescantiae. *J Genet*. 21:207–286.
- Duminil J, Dainou K, Kaviriri DK, Gillet P, Loo J, et al. 2016. Relationships between population density, fine-scale genetic structure, mating system and pollen dispersal in a timber tree from African rainforests. *Heredity*. 116:295–303.
- Field DL, Broadhurst LM, Elliott CP, Young AG. 2017. Population assignment in autopolyploids. *Heredity*. 119:389–401.
- Fisher RA. 1943. Allowance for double reduction in the calculation of genotype frequencies with polysomic inheritance. *Annal Human Genet*. 12:169–171.
- Fisher RA, Mather K. 1943. The inheritance of style length in *Lythrum salicaria*. *Annal Eugen*. 12:1–23.
- Geiringer H. 1949. Chromatid segregation of tetraploids and hexaploids. *Genetics*. 34:665–684.
- Gerber S, Mariette S, Streiff R, Bodenes C, Kremer A. 2000. Comparison of microsatellites and amplified fragment length polymorphism markers for parentage analysis. *Mol Ecol*. 9: 1037–1048.
- Haldane JBS. 1930. Theoretical genetics of autopolyploids. *J Genet*. 22:359–372.
- Huang K, Dunn DW, Ritland K, Li B. 2020. POLYGENE: Population genetics analyses for autopolyploids based on allelic phenotypes. *Methods Ecol Evol*. 11:448–456.
- Huang K, Ritland K, Guo ST, Shattuck M, Li BG. 2014. A pairwise relatedness estimator for polyploids. *Mol Ecol Resour*. 14:734–744.
- Huang K, Wang TC, Dunn DW, Zhang P, Liu RC, et al. 2019. Genotypic frequencies at equilibrium for polysomic inheritance under double-reduction. *G3 (Bethesda)*. 9:1693–1706.
- Huber G. 2016. An investigation of highbush blueberry floral biology and reproductive success in British Columbia, Ph.D. thesis, University of British Columbia.
- Ismail SA, Ghazoul J, Ravikanth G, Kushalappa CG, Shaanker RU, et al. 2017. Evaluating realized seed dispersal across fragmented tropical landscapes: a two-fold approach using parentage analysis and the neighbourhood model. *New Phytol*. 214:1307–1316.
- Kalinowski ST, Taper ML, Marshall TC. 2007. Revising how the computer program CERVUS accommodates genotyping error increases success in paternity assignment. *Mol Ecol*. 16:1099–1106.
- Krebs SL, Hancock JF. 1989. Tetrasomic inheritance of isoenzyme markers in the highbush blueberry. *Vaccin Corymb L Hered*. 63: 11–18.
- Lloyd A, Bomblies K. 2016. Meiosis in autopolyploid and allopolyploid *Arabidopsis*. *Current Opinion Plant Biol*. 30:116–122.
- Manching H, Sengupta S, Hopper KR, Polson SW, Ji Y, et al. 2017. Phased genotyping-by-sequencing enhances analysis of genetic diversity and reveals divergent copy number variants in maize. *G3 (Bethesda)*. 7:2161–2170.
- Marshall TC, Slate J, Kruuk LEB, Pemberton JM. 1998. Statistical confidence for likelihood-based paternity inference in natural populations. *Mol Ecol*. 7:639–655.
- Mather K. 1935. Reductional and equational separation of the chromosomes in bivalents and multivalents. *J Genet*. 30:53–78.
- Monthe FK, Hardy OJ, Doucet J-L, Loo J, Duminil J. 2017. Extensive seed and pollen dispersal and assortative mating in the rain forest tree *Entandrophragma cylindricum* (Meliaceae) inferred from indirect and direct analyses. *Mol Ecol*. 26:5279–5291.
- Muller HJ. 1914. A new mode of segregation in Gregory's tetraploid *Primulas*. *Am Nat*. 48:508–512.
- Neigenfind J, Gyetvai G, Basekow R, Diehl S, Achenbach U, et al. 2008. Haplotype inference from unphased SNP data in heterozygous polyploids based on SAT. *BMC Genomics*. 9:356.
- Norman P, Asfaw A, Tongoona P, Danquah A, Danquah E, et al. 2018. Can parentage analysis facilitate breeding activities in root and tuber crops? *Agriculture*. 8:95.
- Oddou-Muratorio S, Gauzere J, Bontemps A, Rey J-F, Klein EK. 2018. Tree, sex and size: Ecological determinants of male versus female fecundity in three *Fagus sylvatica* stands. *Mol Ecol*. 27:3131–3145.
- Parisod C, Holderegger R, Brochmann C. 2010. Evolutionary consequences of autopolyploidy. *New Phytol*. 186:5–17.
- Pelé A, Rousseau-Gueutin M, Chèvre A-M. 2018. Speciation success of polyploid plants closely relates to the regulation of meiotic recombination. *Front Plant Sci*. 9:907.

- Ravinet M, Westram A, Johannesson K, Butlin R, André C, *et al.* 2016. Shared and nonshared genomic divergence in parallel ecotypes of *Littorina saxatilis* at a local scale. *Mol Ecol.* 25:287–305.
- Ritland K. 2002. Extensions of models for the estimation of mating systems using  $n$  independent loci. *Heredity.* 88:221–228.
- Rodzen JA, Famula TR, May B. 2004. Estimation of parentage and relatedness in the polyploid white sturgeon (*Acipenser transmontanus*) using a dominant marker approach for duplicated microsatellite loci. *Aquaculture.* 232:165–182.
- Schwarz G. 1978. Estimating the dimension of a model. *Ann Statist.* 6:461–464.
- Solís Neffa V, Fernández A. 2000. Chromosome studies in *Turnera* (turneraceae). *Genet Mol Biol.* 23:925–930.
- Stift M, Berenos C, Kuperus P, van Tienderen PH. 2008. Segregation models for disomic, tetrasomic and intermediate inheritance in tetraploids: a general procedure applied to rorippa (yellow cress) microsatellite data. *Genetics.* 179:2113–2123.
- Stift M, Reeve R, Van Tienderen P. 2010. Inheritance in tetraploid yeast revisited: segregation patterns and statistical power under different inheritance models. *J Evol Biol.* 23:1570–1578.
- Taberlet P, Griffin S, Goossens B, Questiau S, Manceau V, *et al.* 1996. Reliable genotyping of samples with very low DNA quantities using PCR. *Nucleic Acids Res.* 24:3189–3194.
- Tambarussi EV, Boshier D, Vencovsky R, Freitas ML, Sebbenn AM. 2015. Paternity analysis reveals significant isolation and near neighbor pollen dispersal in small *Cariniana legalis* Mart. Kuntze populations in the Brazilian Atlantic Forest. *Ecol Evol.* 5: 5588–5600.
- Tan LQ, Liu QL, Zhou B, Yang C-J, Zou X, *et al.* 2019. Paternity analysis using SSR markers reveals that the anthocyanin-rich tea cultivar ‘Ziyan’ is self-compatible. *Sci Horticult.* 245:258–262.
- Wagner AP, Creel S, Kalinowski ST. 2006. Estimating relatedness and relationships using microsatellite loci with null alleles. *Heredity.* 97:336–345.
- Wang J, Scribner KT. 2014. Parentage and sibship inference from markers in polyploids. *Mol Ecol Resour.* 14:541–553.
- Wang JL. 2004. Sibship reconstruction from genetic data with typing errors. *Genetics.* 166:1963–1979.
- Wang JL. 2016. Individual identification from genetic marker data: developments and accuracy comparisons of methods. *Mol Ecol Resour.* 16:163–175.
- Watanabe S, Takakura K-I, Kaneko Y, Noma N, Nishida T. 2018. Skewed male reproductive success and pollen transfer in a small fragmented population of the heterodichogamous tree *Machilus thunbergii*. *J Plant Res.* 131:623–631.
- Weiß CL, Pais M, Cano LM, Kamoun S, Burbano HA. 2018. NQUIRE: a statistical framework for ploidy estimation using next generation sequencing. *BMC Bioinformatics.* 19:122.
- Wood TE, Takebayashi N, Barker MS, Mayrose I, Greenspoon PB, *et al.* 2009. The frequency of polyploid speciation in vascular plants. *Proc Natl Acad Sci USA.* 106:13875–13879.
- Yang H, Chen X, Wong WH. 2011. Completely phased genome sequencing through chromosome sorting. *Proc Natl Acad Sci USA.* 108:12–17.
- Zwart AB, Elliott C, Hopley T, Lovell D, Young A. 2016. Polypatex: an R package for paternity exclusion in autopolyploids. *Mol Ecol Resour.* 16:694–700.

Communicating editor: P. Ingvarsson