



## Accurate genotype imputation in multiparental populations from low-coverage sequence

Zheng, C., Boer, M. P., & van Eeuwijk, F. A.

This is a "Post-Print" accepted manuscript, which has been published in "Genetics"

This version is distributed under a non-commercial no derivatives Creative Commons



([CC-BY-NC-ND](https://creativecommons.org/licenses/by-nc-nd/4.0/)) user license, which permits use, distribution, and reproduction in any medium, provided the original work is properly cited and not used for commercial purposes. Further, the restriction applies that if you remix, transform, or build upon the material, you may not distribute the modified material.

Please cite this publication as follows:

Zheng, C., Boer, M. P., & van Eeuwijk, F. A. (2018). Accurate genotype imputation in multiparental populations from low-coverage sequence. *Genetics*, 210(1), 71-82.  
DOI: 10.1534/genetics.118.300885

You can download the published version at:

<https://doi.org/10.1534/genetics.118.300885>

1 **Accurate genotype imputation in multiparental populations**  
2 **from low-coverage sequence**

3 Chaozhi Zheng\*, Martin P. Boer, and Fred A. van Eeuwijk

4 Biometris, Wageningen University and Research, Wageningen, The Netherlands

5 July 24, 2018

6 **Running head:** Genotype imputation

7 **Key words:** Genotype imputation, genotyping by sequencing, hidden Markov model, cross  
8 pollinated (CP) population, multiparent advanced generation intercross (MAGIC)

9 **\*Corresponding author:**

10 Chaozhi Zheng

11 Biometris

12 Wageningen University and Research

13 PO Box 16, 6700AA Wageningen

14 The Netherlands

15 Email: [chaozhi.zheng@wur.nl](mailto:chaozhi.zheng@wur.nl)

## 16 **Abstract**

17 Many different types of multiparental populations have recently been produced to increase ge-  
18 netic diversity and resolution in quantitative trait loci (QTL) mapping. Low coverage geno-  
19 typing by sequencing (GBS) technology has become a cost effective tool in these populations,  
20 despite large amounts of missing data in offspring and founders. In this work, we present a  
21 general statistical framework for genotype imputation in such experimental crosses from low  
22 coverage GBS data. Generalizing a previously developed hidden Markov model for calculating  
23 ancestral origins of offspring DNA, we present an imputation algorithm that doesn't require  
24 parental data and that is applicable to bi- and multiparental populations. Our imputation al-  
25 gorithm allows heterozygosity of parents and offspring as well as error correction in observed  
26 genotypes. Further, our approach can combine imputation and genotype calling from sequenc-  
27 ing reads, and it also applies to called genotypes from single nucleotide polymorphism (SNP)  
28 array data. We evaluate our imputation algorithm by simulated and real datasets in four dif-  
29 ferent types of populations: the F2, the advanced intercross recombinant inbred lines (AI-RIL),  
30 the multiparent advanced generation intercross (MAGIC), and the cross pollinated (CP) popula-  
31 tion. Because our approach uses marker data and population design information efficiently, the  
32 comparisons with previous approaches show that our imputation is accurate at even very low  
33 ( $< 1\times$ ) sequencing depth, in addition to having accurate genotype phasing and error detection.

## 34 **Introduction**

35 Genotype imputation describes the process of imputing missing genotypes in study individuals,  
36 most often using a high density reference panel of genotypes. For human populations, HapMap  
37 (FRAZER *et al.* 2007) and the 1000 genome project (ALTSHULER *et al.* 2012) provide reference  
38 panels including millions of SNPs. Genotype imputation has become a key step in the genome  
39 wide association studies of human populations to increase the power of QTL detection and  
40 to facilitate meta-analyses of studies at different sets of SNPs (LI and FREUDENBERG 2009;  
41 MARCHINI and HOWIE 2010).

42 Genotype imputation leverages haplotype sharing between study individuals and reference  
43 panels. Along chromosomes, the pattern of haplotype sharing changes due to historical re-  
44 combination. A crucial component of most genotype imputation methods is to infer the local  
45 haplotype clustering and the ancestral haplotypes from reference panels and study individuals  
46 (HOWIE *et al.* 2009; LI *et al.* 2010; BROWNING and BROWNING 2016). The accuracy of im-  
47 putation depends on how well reference panels match study individuals in terms of ancestral  
48 haplotypes (PEI *et al.* 2008; ROSHYARA *et al.* 2016).

49 Next-generation sequencing technology has become an attractive and cost effective tool for  
50 QTL mapping in non-human populations (SPINDEL *et al.* 2013; HEFFELFINGER *et al.* 2014;  
51 KIM *et al.* 2016), and genotype imputation is essential for low coverage sequencing. The focus  
52 of this paper is on experimentally designed populations, particularly for plants, where study  
53 individuals are produced by multi-generation crossing from two or more founders. Many such  
54 multiparental populations have recently been created (e.g. KOVER *et al.* 2009; BANDILLO *et al.*  
55 2013; MACKAY *et al.* 2014; SANNEMANN *et al.* 2015), aiming at increasing genetic diversity  
56 due to many founders and QTL mapping resolution due to accumulated recombination break-  
57 points over multiple generations.

58 The founders of multiparental populations are naturally used as the reference panel for geno-  
59 type imputation. However, there are typically many missing founder genotypes particularly  
60 when both founders and offspring are genotyped by low coverage sequencing, and some of the  
61 founders may even be missing completely (THEPOT *et al.* 2015). In such cases, the population-  
62 based imputation methods (HOWIE *et al.* 2009; LI *et al.* 2010; BROWNING and BROWNING

2016) are not optimal. Alternatively, pedigree-based genotype imputation methods (ABECASIS *et al.* 2002; CHEUNG *et al.* 2013) are computationally intensive if not impossible, because of the large breeding pedigree being often partially or wholly unavailable, and most or all genotypes being missing in intermediate generations.

Recently, several imputation methods were proposed for experimental crosses. XIE *et al.* (2010) described a parent-independent genotyping method for two-way recombinant inbred lines (RILs), where parental genotypes were obtained using a maximum parsimony of recombination. SWARTS *et al.* (2014) described a Full-Sib Family Haplotype Imputation (FSFHap) method for biparental populations, where parental haplotypes were identified by a custom clustering method over non-overlapping windows with a window size of 50 loci along chromosomes. FRAGOSO *et al.* (2016) described a Low-Coverage Biallelic Impute (LB-Impute) algorithm for biparental populations, where parental genotypes were imputed only after offspring genotypes were imputed using a modified Viterbi algorithm over a sliding window (of size 7 loci) along chromosomes. See also HICKEY *et al.* (2015) for genotype imputation in biparental populations in plant breeding.

In experimental crosses, genotype imputation methods have mainly focused on biparental populations. There remain challenges for more complicated experimental designs. HUANG *et al.* (2014) described a genotype imputation method called mpimpute, which is however restricted to the funnel scheme 4- or 8-way RILs. In the funnel scheme, the founders of each line are randomly permuted. In this paper, we present a general statistical framework of genotype imputation from low coverage GBS data, applicable to many scenarios in experimental crosses. First, it applies to both bi- and multiparental populations. Second, it is parent-independent so that it applies even if some founders' genotypes are not available. Third, it integrates with parental phasing and thus applies to mapping populations with outbred founders. Last but not least, it integrates with genotype calling to account for the uncertainties in identifying heterozygous genotypes due to low read numbers.

Our imputation algorithm is called magicImpute, building on a hidden Markov model (HMM) framework that extends our previous work (ZHENG *et al.* 2014; ZHENG 2015; ZHENG *et al.* 2015, 2018). We first evaluate magicImpute with simulated data in four populations: the F2,

92 the AI-RIL, the funnel scheme 8-way RILs, and the CP. Then we analyze four sets of real data:  
 93 the maize F2 (ELSHIRE *et al.* 2011), the maize AI-RIL (HEFFELFINGER *et al.* 2014), the rice  
 94 MAGIC (BANDILLO *et al.* 2013), and the apple CP (GARDNER *et al.* 2014). The term MAGIC  
 95 has been used for many different types of breeding designs, and the rice MAGIC is essentially  
 96 a set of funnel scheme 8-way RILs (BANDILLO *et al.* 2013). In the evaluations by simulation  
 97 and real data, we perform comparisons among magicImpute, Beagle v4.1 (BROWNING and  
 98 BROWNING 2016), LB-impute (FRAGOSO *et al.* 2016) and mpimpute (HUANG *et al.* 2014),  
 99 investigating, among other things, how imputation quality depends on amount of missing data,  
 100 level of homozygosity and coverage of sequencing.

## 101 **Methods**

### 102 **Overview of model**

103 Consider a mapping population derived from a number  $n_F \geq 2$  of founders. We assume that  
 104 linkage groups (chromosomes) are independent, and thus consider only one group. The geno-  
 105 typic data matrix of sampled offspring is denoted by  $\mathbf{y}^O = \{y_{ti}\}_{t=1\dots T, i=1\dots N}$ , with element  $y_{ti}$   
 106 representing the genotype at locus  $t$  in offspring  $i$ . The founder genotype matrix is denoted  
 107 by  $\mathbf{y}^F = \{y_t^F\}_{t=1\dots T}$ , with element  $y_t^F$  being the genotypes at locus  $t$  in all founders. We con-  
 108 sider only bi-allelic markers, and denote the two alleles by 1 and 2. We model either the called  
 109 genotypes from SNP array or GBS data, or the allelic depths of GBS data. The called unphased  
 110 genotype at a locus can take one of six possible values: 11, 12, 22, 1U, 2U, or UU, where U  
 111 denotes an uncertainty allele. For allelic depth data, the genotype is measured by read counts  
 112 for each of two alleles. The ordering and genetic locations of markers are assumed to be known.

113 We build an integrated hidden Markov model for the genotypic data  $\mathbf{y}^O$  and  $\mathbf{y}^F$ , but impute  
 114 missing founder genotypes and missing offspring genotypes separately. The imputation diagram  
 115 and the overview of the HMM are shown in Figure 1. Here the hidden founder haplotype matrix  
 116  $\mathbf{h}^F = \{h_t^F\}_{t=1\dots T}$ , where element  $h_t^F$  is similar to  $y_t^F$  except that it contains information on  
 117 missing genotypes and genotype phases at locus  $t$  in founders. See an example in the following  
 118 section on the genotype model. Conditional on estimated  $\hat{\mathbf{h}}^F$ , the genotypic data for each

119 offspring are analyzed independently by a sub-HMM, with  $x_{ti}$  being the hidden ancestral origin  
 120 state at locus  $t$  in offspring  $i$ . The hidden Markov model will be further explained in the process  
 121 model. See Table 1 for a list of symbols and their brief explanations.

## 122 **The genotype model**

123 **Called genotype:** The genotype model corresponds to the vertical relationships (arrows) in the  
 124 directed acyclic graph of the HMM (Figure 1). Since the genotypes are independent condi-  
 125 tional on the hidden states, we consider a single locus  $t$ . We first model the prior probability  
 126  $P(h_t^F | y_t^F)$ , which is assumed to follow a discrete uniform distribution over all possible com-  
 127 binations under the constraint of called parental genotypes  $y_t^F$ . Consider an example of four  
 128 inbred founders with genotypes at locus  $t$  denoted by 11, 22,  $UU$ , and  $UU$ , respectively. We  
 129 use  $12UU$  as a shorthand for the four homozygous genotypes. Then  $h_t^F$  can take one of four  
 130 possible values 1211, 1212, 1221, 1222 with equal probability. Consider the second example  
 131 of a cross pollinated population, and the genotypes of two outbred parents are denoted by 12  
 132 and  $UU$ . Then  $h_t^F$  can take one of eight possible values 1211, 1212, 1221, 1222, 2111, 2112,  
 133 2121, and 2122, where the last four values account for the alternative phase of the first parent's  
 134 genotype. The founder haplotype matrix  $h^F$  is known if all parental genotypes are observed  
 135 and phased.

136 The hidden founder haplotype  $h_t^F$  is not the true founder haplotype, and it accounts for  
 137 unknown phasing and missing values of called founder genotypes, but not allelic errors. The  
 138 errors in called genotypes can be accounted for in the likelihood  $l_{ti} = P(y_{ti} | h_t^F, x_{ti}, \epsilon_O, \epsilon_F)$   
 139 at locus  $t$  in offspring  $i$ , where  $\epsilon_O$  and  $\epsilon_F$  are the allelic error probabilities for offspring and  
 140 founders, respectively. The calculation of likelihood  $l_{ti}$  has been described in detail in ZHENG  
 141 *et al.* (2015). We describe it briefly as follows. We calculate  $l_{ti}$  by summing over the hidden  
 142 true genotype  $z_{ti}$ , and it holds that

$$l_{ti} = \sum_{z_{ti}} P(y_{ti} | z_{ti}, \epsilon_O) P(z_{ti} | d_{ti}, x_{ti}, \epsilon_F),$$

$$P(z_{ti} | d_{ti}, x_{ti}, \epsilon_F) \propto P(d_{ti} | z_{ti}, x_{ti}, \epsilon_F) P(z_{ti} | x_{ti}),$$

143 where  $d_{ti}$  denotes the derived genotype that is obtained from  $x_{ti}$  and  $h_t^F$  in a deterministic way.  
144 We assign an uninformative prior to  $P(z_{ti}|x_{ti})$ , and calculate  $P(y_{ti}|z_{ti}, \epsilon_O)$  and  $P(d_{ti}|z_{ti}, x_{ti}, \epsilon_F)$ ,  
145 assuming that typing errors occur independently and the observed allele is the alternative one  
146 if an error occurs with probability  $\epsilon_O$  or  $\epsilon_F$ . Here the derived genotype  $d_{ti}$  is the same as true  
147 genotype  $z_{ti}$  if there are no errors in observed founder genotypes ( $\epsilon_F = 0$ ).  
148 **Allelic depth:** We next consider the case that genotypes are represented by allelic depths of  
149 GBS data. We calculate prior probability  $P(h_t^F|y_t^F)$  with  $y_t^F$  being called from founder allelic  
150 depths, where the genotype calling will be described in the next section. For likelihood  $l_{ti}$  at  
151 locus  $t$  in offspring  $i$ , only the calculation of  $P(y_{ti}|z_{ti}, \epsilon_O)$  is different from the case of called  
152 genotypes. We introduce  $\varepsilon$  as the sequencing error probability that is given by  $\varepsilon = 10^{-phred/10}$ ,  
153 where  $phred$  is Phred quality score. The genotype  $y_{ti}$  is represented by  $(r_1, r_2)$ , the number of  
154 reads for alleles 1 and 2, respectively. It holds that

$$\begin{aligned}
P((r_1, r_2)|z' = 11, \varepsilon) &\propto (1 - \varepsilon)^{r_1} \varepsilon^{r_2}, \\
P((r_1, r_2)|z' = 12, \varepsilon) &\propto (1/2)^{r_1+r_2}, \\
P((r_1, r_2)|z' = 21, \varepsilon) &\propto (1/2)^{r_1+r_2}, \\
P((r_1, r_2)|z' = 22, \varepsilon) &\propto \varepsilon^{r_1} (1 - \varepsilon)^{r_2},
\end{aligned} \tag{1}$$

155 conditional on hidden genotype  $z'$  (XIE *et al.* 2010).

156 We interpret  $\epsilon_O$  as a depth-independence allelic error probability, for example, due to the  
157 mis-assignment of reads to the reference genome. And we assume that  $z'$  results from the true  
158 genotype  $z_{ti}$  with error probability  $\epsilon_O$ . Thus,  $P(y_{ti}|z_{ti}, \epsilon_O, \varepsilon)$  can be calculated by summing  
159 over  $z'$  as follows

$$P(y_{ti} = (r_1, r_2)|z_{ti}, \epsilon_O, \varepsilon) = \sum_{z'} P((r_1, r_2)|z', \varepsilon) P(z'|z_{ti}, \epsilon_O)$$

160 where  $P(z'|z_{ti}, \epsilon_O)$  is similar to  $P(y_{ti}|z_{ti}, \epsilon_O)$  in the case of called genotypes, except that  $z'$   
161 is phased. Specifically for  $z_{ti} = 11$ , we have  $P(z'|z_{ti} = 11, \epsilon_O) = (1 - \epsilon_O)^2$ ,  $(1 - \epsilon_O)\epsilon_O$ ,  
162  $\epsilon_O(1 - \epsilon_O)$ , and  $\epsilon_O^2$  for  $z' = 11, 12, 21$ , and  $22$ , respectively. And similarly for  $z_{ti} = 12, 21$ ,  
163 and  $22$ . When there are no ambiguities, we suppress the dependence of  $\varepsilon$  for allelic depth data



164 in the description of the imputation algorithm.

165 **Single genotype calling:** We perform single genotype calling for founder allelic depths of GBS  
166 data before imputation, and for detecting potential erroneous genotypes among offspring during  
167 the last stage of imputation. For single genotype calling from allelic depths, we do not consider  
168 depth-independence errors. The calling is based on the following posterior probability

$$P(z_{ti}|y_{ti} = (r_1, r_2), \varepsilon) \propto P(y_{ti}|z_{ti}, \varepsilon)P(z_{ti}),$$

169 where  $P(y_{ti}|z_{ti}, \varepsilon)$  is given by Equation 1 and  $P(z_{ti}) = 1/4, 1/2,$  and  $1/4$  for  $z_{ti} = 11, 12,$   
170 and  $22,$  respectively. Note that  $z_{ti}$  is unphased only in case of single genotype calling, and it is  
171 phased elsewhere. The genotype with posterior probability being greater than threshold  $P_{call}$  is  
172 called. If no genotype is called, we calculate the posterior probability

$$P(z_{ti} = 1U|y_{ti}, \varepsilon) = P(z_{ti} = 11|y_{ti}, \varepsilon) + P(z_{ti} = 12|y_{ti}, \varepsilon),$$

$$P(z_{ti} = 2U|y_{ti}, \varepsilon) = P(z_{ti} = 22|y_{ti}, \varepsilon) + P(z_{ti} = 12|y_{ti}, \varepsilon).$$

173 The genotype  $1U$  is called if  $P(z_{ti} = 1U|y_{ti}, \varepsilon) > P_{call}$  and  $P(z_{ti} = 1U|y_{ti}, \varepsilon) > P(z_{ti} =$   
174  $2U|y_{ti}, \varepsilon),$  and similarly for genotype  $2U$ . The genotype is set to  $UU$  if no calling occurs.

## 175 **The process model**

176 The process model corresponds to the horizontal relationships (arrows) in the directed acyclic  
177 graph of the HMM (Figure 1). It has been described in detail (ZHENG *et al.* 2014; ZHENG  
178 2015; ZHENG *et al.* 2015), and we give a brief summary in the following. The process  $\{x_{ti}\}_{t=1}^T$   
179 for offspring  $i$  describes how the ancestral origins change along chromosomes. At a locus  $t,$   
180 let  $x_{ti} = (x_{ti}^m, x_{ti}^p)$  be the ancestral origins on the maternally (m) and paternally (p) derived  
181 chromosomes. If offspring  $i$  is fully inbred, we have  $x_{ti}^m = x_{ti}^p$  so that the ancestral origin  
182 process along the maternally derived chromosome is the same as the process along the pater-  
183 nally derived chromosome, and it is thus termed "depModel". On the other hand, if offspring  
184  $i$  is completely outbred, the ancestral origin process along the maternally derived chromosome  
185  $\{x_{ti}^m\}_{t=1}^T$  is independent of the process  $\{x_{ti}^p\}_{t=1}^T$  along the paternally derived chromosome, and

186 it is therefore termed "indepModel". In the general model called "jointModel",  $x_{ti}^m$  and  $x_{ti}^p$  are  
 187 modeled jointly. We have kept the model terms (e.g. "jointModel") consistent with ZHENG  
 188 *et al.* (2015).

189 In all three models, the ancestral origin process along two chromosomes is assumed to  
 190 follow a Markov process, so that the ancestral origins  $x_{ti}$  at locus  $t$  depends only on  $x_{t-1,i}$  at  
 191 locus  $t - 1$  but not on the previous  $\{x_{t',i}\}_{t'=1}^{t-2}$ . Thus, the joint prior distribution of  $\{x_{ti}\}_{t=1}^T$  can  
 192 be specified by the initial distribution  $\pi(x_{1i})$  and the transition probability  $P(x_{ti}|x_{t-1,i})$  at  $t =$   
 193  $2, \dots, T$ . The initial distribution  $\pi(x_{1i})$  is specified by the stationary distribution of the Markov  
 194 process, so that the prior process model does not depend on the direction of chromosomes.  
 195 The initial distribution  $\pi(x_{1i})$  and transition probability  $P(x_{ti}|x_{t-1,i})$  can be specified from the  
 196 breeding design of a mapping population, that is, how the sampled offspring is produced from  
 197 the founders; the transition probability also depends on inter-marker distances. See ZHENG  
 198 *et al.* (2014), ZHENG (2015), and ZHENG *et al.* (2018) for the details of calculating  $\pi(x_{1i})$  and  
 199  $P(x_{ti}|x_{t-1,i})$  under various breeding designs.

## 200 **Founder imputation**

201 Because the state space of the HMM exponentially increases with the number  $N$  of sampled  
 202 offspring, the exact inference of the founder haplotype matrix  $\mathbf{h}^F$  is computationally intractable,  
 203 even using the forward-backward algorithm (RABINER 1989). In the following, we describe  
 204 an approximate forward-backward procedure for maximum likelihood estimation of  $\mathbf{h}^F$ . Our  
 205 forward algorithm calculates recursively the posterior probabilities  $\gamma(h_t^F)$  and  $\alpha(x_{ti}|h_t^F)$  for  
 206 offspring  $i = 1, \dots, N$ , conditional on genotypic data up to locus  $t$ . It proceeds as follows:

207 **A0 Initialize at  $t = 1$**

$$\begin{aligned}\tilde{\alpha}(x_{1i}|h_1^F) &= P(y_{1i}|h_1^F, x_{1i}, \epsilon_O, \epsilon_F)\pi(x_{1i}), \\ \gamma(h_1^F) &\propto P(h_1^F|y_1^F) \prod_{i=1}^N \sum_{x_{1i}} \tilde{\alpha}(x_{1i}|h_1^F), \\ \alpha(x_{1i}|h_1^F) &= \tilde{\alpha}(x_{1i}|h_1^F) / \sum_{x_{1i}} \tilde{\alpha}(x_{1i}|h_1^F).\end{aligned}$$

208 A1 For  $t = 2, \dots, T$

$$\begin{aligned}\tilde{\alpha}(x_{ti}|h_t^F) &= P(y_{ti}|h_t^F, x_{ti}, \epsilon_O, \epsilon_F) \sum_{x_{t-1,i}} P(x_{ti}|x_{t-1,i}) \sum_{h_{t-1}^F} \gamma(h_{t-1}^F) \alpha(x_{t-1,i}|h_{t-1}^F), \\ \gamma(h_t^F) &\propto P(h_t^F|y_t^F) \prod_{i=1}^N \sum_{x_{ti}} \tilde{\alpha}(x_{ti}|h_t^F), \\ \alpha(x_{ti}|h_t^F) &= \tilde{\alpha}(x_{ti}|h_t^F) / \sum_{x_{ti}} \tilde{\alpha}(x_{ti}|h_t^F),\end{aligned}$$

209 where  $\tilde{\alpha}(x_{ti}|h_t^F)$  is an unnormalized probability, and the normalization constant for  $\gamma(h_t^F)$  is  
 210 not shown. The key approximation comes from the independence of offspring in the calcula-  
 211 tion of  $\gamma(h_t^F)$ . ZHENG *et al.* (2016) have described a similar forward algorithm for haplotype  
 212 reconstruction in tetraploid populations.

213 The maximum likelihood estimation of founder haplotypes is based on the posterior prob-  
 214 abilities  $\alpha(x_{ti}|h_t^F)$  and  $\gamma(h_t^F)$  from algorithm A. The maximization proceeds backwardly as  
 215 follows:

216 B0 Initialize at  $t = T$ :  $\hat{h}_T^F = \operatorname{argmax} \gamma(h_T^F)$  and  $\hat{x}_{T,i} = \operatorname{argmax} \alpha(x_{T,i}|h_T^F)$  for  $i = 1, \dots, N$ .

217 B1 For  $t = T - 1, \dots, 1$

$$\begin{aligned}\beta(x_{ti}|h_t^F) &= \alpha(x_{ti}|h_t^F) P(\hat{x}_{t+1,i}|x_{ti}), \\ \hat{h}_t^F &= \operatorname{argmax} \gamma(h_t^F) \prod_{i=1}^N \sum_{x_{ti}} \beta(x_{ti}|h_t^F), \\ \hat{x}_{ti} &= \operatorname{argmax} \beta(x_{ti}|\hat{h}_t^F).\end{aligned}$$

218 It is possible that multiple argument values correspond to the same maximum. If such ties  
 219 occur, we randomly choose one of these values. FRIEL and RUE (2007) have described a  
 220 similar backward maximization algorithm for general factorisable models.

221 Preliminary simulations showed that our forward-backward procedure is occasionally less  
 222 accurate on the left end of chromosomes in case of sparse data. We overcome this problem  
 223 by two rounds of maximization. Specifically, we fix the founder haplotypes on the right-half  
 224 chromosomes ( $t > T/2$ ) after the first round of maximization, and then perform the second

225 round with reversed chromosome direction.

## 226 **Offspring imputation**

227 Conditional on the imputed founder haplotype matrix  $\hat{\mathbf{h}}^F$ , all the offspring are independent.  
228 For each offspring, we first perform the posterior decoding algorithm to calculate the posterior  
229 probabilities of ancestral origins at all loci (RABINER 1989; ZHENG *et al.* 2015). Then we  
230 calculate the posterior probabilities of true genotypes, from which missing genotypes can be  
231 imputed.

232 We obtain  $P(z_{ti}|\mathbf{y}^O, \hat{\mathbf{h}}^F, \epsilon_O, \epsilon_F)$  by marginalizing the following joint posterior probability

$$P(z_{ti}, x_{ti}|\mathbf{y}^O, \hat{\mathbf{h}}^F, \epsilon_O, \epsilon_F) = P(z_{ti}|d_{ti}, x_{ti}, \epsilon_F)P(x_{ti}|\mathbf{y}^O, \hat{\mathbf{h}}^F, \epsilon_O, \epsilon_F),$$

233 where the posterior probability  $P(x_{ti}|\mathbf{y}^O, \hat{\mathbf{h}}^F, \epsilon_O, \epsilon_F)$  can be calculated by the function magi-  
234 cReconstruct in the RABBIT software (ZHENG *et al.* 2015), which has been extended to analyze  
235 allelic depths of GBS data. Here the derived genotype  $d_{ti}$  is completely determined by  $x_{ti}$  and  
236  $\hat{h}_t^F$ , and the calculation of  $P(z_{ti}|d_{ti}, x_{ti}, \epsilon_F)$  has been described in the genotype model.

237 From the marginal posterior probability  $P(z_{ti}|\mathbf{y}^O, \hat{\mathbf{h}}^F, \epsilon_O, \epsilon_F)$ , we perform both imputation  
238 and error detection for offspring  $i$ . For imputation, the missing genotype in offspring  $i$  at lo-  
239 cus  $t$  is imputed to be  $\hat{z}_{ti}$  if its marginal posterior probability is larger than a given threshold  
240  $P_{impute}$ . For error detection, the observed called genotype  $y_{ti}$  is corrected if the most probable  
241 genotype is different from  $y_{ti}$  and the maximal marginal posterior probability is larger than a  
242 given threshold  $P_{detect}$ .

## 243 **Data simulation**

244 We simulate sequence data, mimicking real data in the following mapping populations: the AI-  
245 RIL, the F2, the MAGIC(funnel scheme 8-way RIL), and the CP. These populations differ in  
246 the number of founders and the heterozygosity level of founders and offspring (Table 2). For  
247 each type of mapping population, we simulate independently three sample sizes: 100, 200, and  
248 500, that is, the number of sampled offspring in the last generation. Independently for each type

249 of population with a given sample size, we first simulate the breeding pedigree according to the  
250 corresponding real data. The AI-RIL consists of five generations of random mating starting  
251 from the F1 generation and six generations of selfing; the size of the random mating population  
252 is set to 1000. For each offspring of the MAGIC, the founders are randomly permuted so that  
253 the number of funnels equals the sample size.

254 Given a breeding pedigree for each mapping population, we assign a unique founder genome  
255 label (FGL) to each inbred founder or to the haploid gamete of each outbred founder. We  
256 simulate only one linkage group. Each offspring gamete is a random mosaic of FGL blocks  
257 determined by chromosomal crossovers between two parental chromosomes. The number of  
258 crossovers in a gamete follows a Poisson distribution with mean being the chromosome length  
259 in Morgan, and the positions of crossovers are uniformly distributed across the chromosome.  
260 We set true founder haplotypes based on the founders imputed from the available real data  
261 (see Table 2), and obtain the true offspring genotypes by replacing FGLs with the true founder  
262 haplotypes. We apply the same error model to the true founder haplotypes with  $\epsilon_F = 0.005$  and  
263 to the true offspring genotypes with  $\epsilon_O = 0.005$ .

264 We simulate read count data for each obtained founder or offspring genotype. Independently  
265 for each allele of a genotype, the number of reads is assumed to follow an exponential  
266 distribution with mean being  $\lambda/2$ , where we set  $\lambda = 8$ ; the number of erroneous reads follow  
267 a binomial distribution with probability  $\epsilon = 0.001$ , and the erroneous read corresponds to the  
268 alternative allele. The allelic depths of genotypes are obtained by combining reads of the two  
269 alleles. The allelic depths of founder and offspring genotypes are re-set to be missing with prob-  
270 abilities 0.25 and 0.15, respectively. We obtain 12 full datasets, 3 population sizes for each of  
271 the four mapping populations, with average offspring read depth 6.8. To study the dependence  
272 of sequencing coverage, we retain the same founder reads and randomly sample offspring reads  
273 with probability  $2^{-i}$  for  $i = 0, 1, \dots, 10$ , resulting in a total of 132 test datasets.

## 274 **Real data**

275 Table 2 shows a summary of real data after filtering. For the maize AI-RIL (HEFFELFINGER  
276 *et al.* 2014) and the maize F2 (ELSHIRE *et al.* 2011), we use the GBS data that have been

277 prepared by FRAGOSO *et al.* (2016) as the input data of LB-Impute. For the rice MAGIC  
278 (BANDILLO *et al.* 2013), we use the called genotypes that have been prepared by HUANG *et al.*  
279 (2014) for mpimpute. For the apple CP (GARDNER *et al.* 2014), we filter the original allelic  
280 depth data by removing markers with the missing fraction of called genotypes larger than 50%,  
281 and removing markers with segregation distortion at significant level 0.01. During the filtering  
282 process, a single genotype is called with threshold  $P_{call} = 0.99$  and 0.95 for founders and  
283 offspring, respectively, as described in the previous section on single genotype calling. And the  
284 quality score is set to  $phred = 30$  so that the sequencing error probability  $\varepsilon = 10^{-phred/10} =$   
285 0.001.

286 To calculate imputation accuracy, we mask a subset of high-confidence genotypes and use  
287 them as the pseudo-true genotypes. For the GBS data, the genotypes are first called with a very  
288 large threshold  $P_{call} = 0.9999$  and the quality scores being 30 and 40 for apple and maize,  
289 respectively. The called genotypes (excluding  $UU$ ,  $1U$  and  $2U$ ) are masked with probability  
290 being 0.25 and 0.05 for founders and offspring, respectively. After masking, the fractions of  
291 founder genotypes without reads are 0.23, 0.24, and 0.19 for the maize AI-RIL, the maize F2,  
292 and the apple CP, respectively. And the fractions of offspring genotypes without reads are 0.77,  
293 0.16, and 0.095. For each of three masked full datasets, we retain the same founder reads  
294 and randomly sample offspring reads with probability  $2^{-i}$  for  $i = 0, 1, \dots, 10$ , resulting in 33  
295 real sequencing datasets. For the called genotypes of the rice MAGIC, the missing fraction of  
296 founder genotypes after masking is 0.3. From this masked dataset, five datasets are produced  
297 independently by masking called offspring genotypes to give missing fractions from 0.5 to 0.9  
298 at step size 0.1.

## 299 **Algorithm evaluation**

300 To set up the algorithm magicImpute, we perform sensitivity analysis of  $P_{impute}$ ,  $P_{detect}$ , and  
301  $\varepsilon_O$ . For each mapping population with size 200 and read depth 0.85, we impute the simulated  
302 dataset with the input data being called genotypes and the first two founders' genotypes being  
303 not available. By default, we set  $\varepsilon_F = 0.005$ , and the input genotypes are called from allelic  
304 depths with threshold  $P_{call} = 0.99$  and 0.95 for founders and offspring, respectively. Figures

305 S1&S2 show that the accuracies of imputation and error detection increase slightly with  $P_{impute}$   
306 from 0.6 to 0.95, while the fractions of imputation and error detection decrease slightly. Figures  
307 S1&S2 also show that the performances of imputation and error detection often become a bit  
308 worse when  $\varepsilon_O$  increases by a factor of 10. The effects of these parameters are marginal in  
309 general. Thus we set somewhat arbitrarily  $P_{impute} = 0.9$ ,  $P_{detect} = 0.9$ , and  $\varepsilon_O = 0.005$  in the  
310 following evaluations. The algorithm `magicImpute` also outputs the posterior probabilities of  
311 all possible genotypes for all offspring at all markers, from which we can perform imputation  
312 and error detection with different  $P_{impute}$  and  $P_{detect}$ .

313 We evaluate `magicImpute` by both simulated and real data in the four types of mapping  
314 populations. For each of the simulated datasets and the real GBS datasets, we run `magicImpute`  
315 in the four combinations: the first two founders' genotypes are available or not, and the input  
316 data are allelic depths or called genotypes. Here the quality scores are 30 for the simulated data  
317 and the real maize GBS data, and 40 for the real apple GBS data. For the real rice data, we  
318 run `magicImpute` in the two combinations: the first two founders' genotypes are available or  
319 not. Results of `magicImpute` are compared with those of Beagle v4.1 in all populations. We  
320 run Beagle v4.1 for the called genotypes in two ways: without reference panels and use the  
321 founder haplotypes imputed by `magicImpute` as the reference panels. Additionally, we run LB-  
322 Impute for the biparental populations AI-RIL and F2 with the input data being allelic depths,  
323 and run `mpimpute` for the MAGIC population with the input data being called genotypes. LB-  
324 Impute and `mpimpute` do not work if some founders' genotypes are not available. The running  
325 settings of `magicImpute`, Beagle v4.1, LB-Impute, and `mpimpute` are described in Supporting  
326 Information, File S1. See SWARTS *et al.* (2014) and FRAGOSO *et al.* (2016) for comparisons of  
327 FSFHap with Beagle and LB-Impute.

## 328 **Data availability**

329 The algorithm `magicImpute` is implemented in Mathematica 11.0 (WOLFRAM RESEARCH  
330 2016), and it has been included as a function in the RABBIT software. RABBIT is available  
331 at <https://github.com/chaozhi/RABBIT.git>, and it is offered under the GNU Af-  
332 fero general public license, version 3 (AGPL-3.0). Example scripts for simulating genotypic

333 data are included. The real maize AI-RIL and F2 data have been described by HEFFELFINGER  
334 *et al.* (2014) and ELSHIRE *et al.* (2011), respectively, and they have been prepared by FRAGOSO  
335 *et al.* (2016) for LB-Impute. The rice MAGIC data have been described by BANDILLO *et al.*  
336 (2013), and they have been prepared by HUANG *et al.* (2014) for mpimpute. The apple CP data  
337 are available from GARDNER *et al.* (2014).

## 338 **Results**

### 339 **Simulation evaluation**

340 Figures 2-4 and Figures S3-S7 show the comparisons among magicImpute, Beagle, LB-Impute,  
341 and mpimpute in terms of imputation accuracy, error detection, and genotype phasing. All  
342 results are obtained from the simulated populations of size 200, except Figure S4 that shows the  
343 effects of population size.

344 **Imputation accuracy:** Figures 2 and S3 show the comparisons of imputation accuracy. One  
345 of the most striking patterns is that there exist break points for magicImpute and Beagle but not  
346 for LB-Impute and mpimpute. As shown in Figure 2 for the imputation accuracy of offspring  
347 genotypes, the break points of magicImpute are 0.053, 0.11, 0.21, and 0.21 read depth for the  
348 AI-RIL, the F2, the MAGIC, and the CP, respectively, much lower than the break points of 0.42,  
349 3.4, 0.85, and 3.4 read depth for Beagle. As shown in the left panels of Figure S3, the break  
350 points of magicImpute for founder imputation are the same as those for offspring imputation;  
351 Beagle does not impute founder genotypes.

352 As for mpimpute and LB-Impute, they perform slightly worse than magicImpute. The im-  
353 putation accuracy of mpimpute is  $\sim 1.7\%$  lower than that of magicImpute when read depth  $>$   
354 0.21 (Figure 2C). The imputation accuracies of LB-Impute at the highest read depth are similar  
355 to those of magicImpute, but they decrease gradually with decreasing read depth. In addition,  
356 the imputation fractions of LB-Impute at the highest read depth are around 0.8, much smaller  
357 than those of magicImpute (Figure S3B&D).

358 The unavailability of the first two founders' genotypes has no noticeable effects on the  
359 performance of magicImpute for the AI-RIL, the F2, and the MAGIC, as long as read depth is



360 higher than the break point. However for the CP, the availability of the two outbred founders'  
361 genotypes results in  $\sim 2\%$  lower accuracy of imputing founder genotypes (Figure S3G), due  
362 to the calling errors in the available founder genotypes. As a result, the imputation accuracy of  
363 offspring genotypes is  $\sim 4\%$  lower (Figure 2D).

364 Whether the input data are allelic depths or called genotypes has little influence on the  
365 performance of magicImpute. However for the almost homozygous populations AI-RIL and  
366 MAGIC, the ceiling limit of imputation accuracy decreases with increasing read depth instead  
367 of leveling off (Figure 2A&C). This is due to the assumption of homozygosity during the prior  
368 genotype calling, and the information on residual heterozygosity is lost after transforming al-  
369 lelic depths into called genotypes. The percentage of heterozygotes among missing genotypes  
370 increases with increasing read depth, and they are always missing and wrongly imputed.

371 Figure S4 shows that the main effect of population size is shifting the break points of the  
372 imputation accuracy obtained by magicImpute and Beagle.

373 **Error detection:** We evaluate the error detection of magicImpute in the case of the input data  
374 being called genotypes. A suspicious genotype error is detected by magicImpute when the most  
375 probable true genotype is different from the input called genotype and the maximum posterior  
376 probability is larger than the default threshold  $P_{detect} = 0.9$ . As shown in Figures 3 and S5,  
377 the unavailability of the first two founders' genotypes greatly improve the error detections for  
378 the F2, the CP, and the AI-RIL, but it has little effects on the MAGIC with multiple founders.  
379 This indicates that the errors in the available founder genotypes adversely affect the detection  
380 of offspring genotypes.

381 Figures 3 and S5 show that the error detection in the almost homozygous populations AI-  
382 RIL and the MAGIC is much worse than in the F2 and the CP. This is due to the homozygosity  
383 assumption under which the input genotypes are being called for the AI-RIL and the MAGIC;  
384 most offspring genotype errors are heterozygous and they cannot be detected and corrected  
385 when the heterozygosity information is lost during the prior genotype calling. Figure S6 shows  
386 that the error detection in the AI-RIL and the MAGIC is much better when homozygosity is not  
387 assumed.

388 **Genotype phasing:** We evaluate the phasing accuracy for the heterozygous populations F2 and

389 CP obtained by magicImpute and Beagle; mpimpute and LB-impute do not perform phasing.  
390 The phasing accuracy is measured in two ways: the switch accuracy is defined as one minus the  
391 number of switches divided by the number of opportunities for switch error, and the heterozy-  
392 gous accuracy denotes the percentage of correctly phased heterozygous genotypes. A switch  
393 error occurs if the heterozygous genotype at a site has phase switched related to that of the  
394 previous heterozygous site.

395 As shown in Figures 4 and S7, the phasing accuracy has similar patterns and the same break  
396 points as those of the imputation accuracy (Figure 2) for magicImpute and Beagle, so that the  
397 phasing of magicImpute is more robust to missing data. For the CP, the switch accuracy and the  
398 heterozygous accuracy of magicImpute are close to 1 when read depth is higher than the break  
399 point, whereas the heterozygous accuracy of Beagle is less than 0.8. The difference between  
400 switch and heterozygous accuracy indicates that the wrongly phased heterozygous genotypes  
401 occur in blocks and they could be corrected by a few switches between the two haplotypes  
402 within an offspring.

403 Figures 4 and S7 show that the availability of the two founders' genotypes are unimportant  
404 to genotype phasing. The phasing accuracy of Beagle increases slightly when read depth is  
405 higher than the break point. However for magicImpute in the CP, the ceiling limit of phasing  
406 accuracy decreases a bit, consistent with the decrease of ceiling imputation accuracy because of  
407 the errors in the available founder genotypes.

## 408 **Evaluation by real data**

409 Figures 5 and S8 show the results of genotype imputation obtained from the real data in the  
410 four mapping populations. Error detection and genotype phasing cannot be evaluated since true  
411 genotypes and phases are not available; the imputation accuracy is calculated based on masked  
412 genotypes. Figure 5 shows the patterns similar to those of the simulation evaluation. The break  
413 points for magicImpute are at much lower read depths or larger missing fractions than those of  
414 Beagle. The magicImpute accuracy is slightly larger than that of mpimpute, and it is always  
415 high until the break point. In contrast to that, the LB-Impute accuracy decreases gradually with  
416 read depth.

417 **Maize AI-RIL and F2:** Figure 5A&B and Figure S8A-D show the results of genotype im-  
418 putation in the real biparental populations AI-RIL and F2. For magicImpute, the offspring  
419 imputation accuracies at the highest read depth are higher than 0.980 in the AI-RIL and 0.987  
420 in the F2. The corresponding accuracies are 0.970 and 0.986 for Beagle, whereas they are 0.917  
421 and 0.986 for LB-Impute. The imputation fractions at the highest read depth for both magicIm-  
422 pute and Beagle are larger than 0.960, whereas for LB-Impute they are 0.720 in the AI-RIL and  
423 0.906 in the F2.

424 FRAGOSO *et al.* (2016) obtained the imputation accuracies 0.970 for the AI-RIL and 0.946  
425 for the F2, and the differences may be due to the masking of founder genotypes and the usage  
426 of a small genotype error probability for magicImpute.

427 **Rice MAGIC:** Figure 5C shows that the imputation accuracies of magicImpute and mpimpute  
428 are almost independent of missing fraction of the input offspring genotypes in the range from  
429 0.5 to 0.9. On average, the offspring imputation accuracy of magicImpute is higher than that  
430 of mpimpute by 2.5%. The Beagle imputation accuracy is comparable to that of magicImpute  
431 when the missing fraction is no greater than the break point of 0.7.

432 Figure S8E shows that the founder imputation accuracies are around 0.94 and 0.89 for  
433 mpimpute and magicImpute, respectively, whereas they are close to 1 in the simulation eval-  
434 uation. The imputation fraction of founder genotypes for mpimpute gradually decreases from  
435 0.947 to 0.922 with increasing missing fraction (Figure S8E); magicImpute imputes all missing  
436 founder genotypes. As a result, the offspring imputation fraction of mpimpute decreases rapidly  
437 from 0.92 to 0.6, whereas it is always around 0.96 for magicImpute (Figure S8F).

438 **Apple CP:** Figure 5D shows the results of offspring imputation accuracy obtained from the real  
439 apple data. The imputation accuracy of magicImpute decreases from 0.94 to 0.88 when read  
440 depth decreases from 15 to 0.46, in comparison with the almost constant accuracy of 0.96 in  
441 the simulated results in Figure 2D. The Beagle imputation accuracy is comparable to that of  
442 magicImpute, when read depth is no less than the break point of 3.7.

443 As shown in Figure S8G, the founder imputation accuracy of magicImpute at the highest  
444 read depth is around 0.96 when the two founders' genotypes are available, whereas it decreases  
445 to 0.75 when the two founders' genotypes are missing. The low accuracy is very likely because

446 of the mix up of the imputed genotypes between the two founders.

447 **Running time:** The running times for the four real datasets at the highest read depths or the  
448 smallest missing fractions are given in Table 1. Beagle is fastest in all populations. For the  
449 biparental populations, LB-Impute is much slower than magicImpute. And for the rice MAGIC,  
450 mpimpute is similar to Beagle, and faster than magicImpute.

451 The main computational load of magicImpute is the first two steps for founder imputation  
452 and phasing (Figure 1). The founder imputation of mpimpute and LB-impute is based on the  
453 decoding algorithm of the sub-HMM for each offspring, corresponding to the third step of  
454 magicImpute.

## 455 Discussion

456 We have implemented an HMM framework magicImpute for genotype imputation from low  
457 coverage sequence or SNP array data. The evaluations by simulation and real data in the four  
458 types of mapping populations demonstrate that magicImpute is accurate and flexible, despite  
459 the population being multiparental, founders being missing, founders being heterozygous, off-  
460 spring being heterozygous, or sequencing coverage being low. The simulation evaluations also  
461 demonstrate the good performance of magicImpute for error detection and genotype phasing.

462 Although the dependence of imputation accuracy on sequence coverage varies with popu-  
463 lation size, marker density, and distribution of reads, magicImpute performs much better than  
464 Beagle, LB-Impute, and mpimpute at very low coverage. Beagle breaks down at much higher  
465 read depth in heterozygous populations than in almost homozygous populations, probably be-  
466 cause of unsuccessful pre-phasing of Beagle imputation for heterozygous populations. Alter-  
467 native pre-phasing methods might increase the follow-up imputation accuracy (WHALEN *et al.*  
468 2017). The LB-Impute accuracy in biparental populations decreases with decreasing read depth,  
469 probably because the number of markers in the Markov trellis window is only 7 by default (large  
470 window size would result in dramatic increases in running time). The lower LB-Impute accu-  
471 racy in the real AI-RIL than in the simulated AI-RIL may be due to the heavy tailed distribution  
472 of read depth in the real data and its inability of borrowing distant marker information.

473 Low coverage sequencing can be represented as allelic depths or called genotypes for the

474 input of magicImpute. The simulation and real evaluations show that the prior transformation of  
475 allelic depths into called genotypes has no appreciable effects, if homozygosity is not assumed  
476 for the transformation in almost homozygous populations. It indicates that little information  
477 is lost in the prior transformation, where the two half called genotypes ( $1U$  and  $2U$ ) keep se-  
478 quence read information efficiently. Genotype likelihoods, a probabilistic representation of low  
479 coverage sequencing, have been alternatively used in many imputation methods such as Beagle  
480 v4.1.

481 It is implicitly assumed by magicImpute that sequencing reads are too short to cover more  
482 than two polymorphic sites, and the phasing information of long reads is ignored. Thus magicIm-  
483 pute would not rely on long reads. For very low coverage sequencing, the distances between  
484 detected neighbor polymorphic sites are expected to be too long, and very long reads are thus  
485 required to keep the phasing information. On the other hand, our HMM imputation framework  
486 provides a solid step for the extension to utilize phasing information.

487 One key assumption of magicImpute is no segregation distortion, when incorporating breed-  
488 ing design information into the HMM. The assumption is not expected to be a problem for bi-  
489 parental populations with only two inbred founders, as confirmed in our real data evaluation.  
490 For the MAGIC and the CP, the founder imputation accuracies in the real data evaluations are  
491 lower than simulation results, probably because of segregation distortion in the real data. For  
492 real MAGIC, magicImpute has higher offspring imputation accuracy and lower founder impu-  
493 tation accuracy than mpimpute, indicating that the offspring imputation is not affected by the  
494 possible segregation distortion.

495 Secondly, magicImpute assumes that the input genetic map is correct, as do Beagle, LB-  
496 Impute, and mpimpute. The assumption contributes to the differences of ceiling offspring impu-  
497 tation accuracy between simulation and real data evaluations. For the real apple CP, GARDNER  
498 *et al.* (2014) estimated the proportion of markers that are inconsistent with the physical grouping  
499 is as high as 18.3%, which might explain why the accuracy is relatively low (from 0.88 to 0.94)  
500 when read depth is no less than the break point (Figure 5D). See for example MONEY *et al.*  
501 (2015) and RUTKOSKI *et al.* (2013) for map-independent imputations in association panels.

502 Another assumption of magicImpute is on the conditional independence of offspring. In the

503 approximate forward algorithm for founder imputation, offspring are assumed to be independent  
504 given the posterior probabilities up to the current time. This approximation is well validated by  
505 the very accurate founder imputation in the simulation evaluations. Conditional on the imputed  
506 founder haplotypes, offspring are assumed to be independent, which is not always true because  
507 these offspring share parents in the intermediate generations. The algorithm magicImpute partly  
508 accounts for this relationship by the pre-calculated HMM parameters based on available breed-  
509 ing pedigrees, and thus the offspring imputation utilizes the marker information of the others  
510 indirectly via the founder imputation.

511 In conclusion, we have demonstrated that magicImpute is more accurate and robust to low  
512 sequencing depth than the current methods, because magicImpute can incorporate experimental  
513 design and utilize marker data efficiently. Furthermore, magicImpute is not restricted to specific  
514 experimental designs, and it can perform parental imputation and phasing in situations where  
515 most current methods are incapable.

## 516 **Acknowledgments**

517 The authors thank Emma Huang for helps on mpimpute, Cris Wijnen for valuable discussion  
518 on sequencing technology, and three anonymous reviewers for their constructive comments.  
519 This research was supported by the Stichting Technische Wetenschappen (STW) - Technology  
520 Foundation, which is part of the Nederlandse Organisatie voor Wetenschappelijk Onderzoek -  
521 Netherlands Organisation for Scientific Research, and which is partly funded by the Ministry of  
522 Economic Affairs. The specific grant number was STW-Rijk Zwaan project 12425.

## 523 **Author contributions**

524 CZ designed the study, created the model, developed the software and algorithm, and wrote  
525 the first draft of the manuscript. MPB and FAE provided critical feedback, helped shape the  
526 manuscript, and acquitted financial support. FAE supervised the project. All authors read and  
527 approved the final manuscript.

## Literature Cited

- 528
- 529 ABECASIS, G. R., S. S. CHERNY, W. O. COOKSON, and L. R. CARDON, 2002 Merlin-rapid  
530 analysis of dense genetic maps using sparse gene flow trees. *Nature Genetics* **30**: 97–101.
- 531 ALTSHULER, D. M., R. M. DURBIN, G. R. ABECASIS, D. R. BENTLEY, A. CHAKRAVARTI,  
532 *et al.*, 2012 An integrated map of genetic variation from 1,092 human genomes. *Nature* **491**:  
533 56–65.
- 534 BANDILLO, N., C. RAGHAVAN, P. A. MUYCO, M. A. L. SEVILLA, I. T. LOBINA, *et al.*,  
535 2013 Multi-parent advanced generation inter-cross (MAGIC) populations in rice: progress  
536 and potential for genetics research and breeding. *RICE* **6**: 11.
- 537 BROWNING, B. L., and S. R. BROWNING, 2016 Genotype imputation with millions of refer-  
538 ence samples. *American Journal of Human Genetics* **98**: 116–126.
- 539 CHEUNG, C. Y. K., E. A. THOMPSON, and E. M. WIJSMAN, 2013 GIGI: an approach to  
540 effective imputation of dense genotypes on large pedigrees. *American Journal of Human*  
541 *Genetics* **92**: 504–516.
- 542 ELSHIRE, R. J., J. C. GLAUBITZ, Q. SUN, J. A. POLAND, K. KAWAMOTO, *et al.*, 2011 A  
543 robust, simple genotyping-by-sequencing (GBS) approach for high diversity species. *PLOS*  
544 *ONE* **6**: e19379.
- 545 FRAGOSO, C. A., C. HEFFELFINGER, H. Y. ZHAO, and S. L. DELLAPORTA, 2016 Imputing  
546 genotypes in biallelic populations from low-coverage sequence data. *Genetics* **202**: 487–495.
- 547 FRAZER, K. A., D. G. BALLINGER, D. R. COX, D. A. HINDS, L. L. STUVE, *et al.*, 2007 A  
548 second generation human haplotype map of over 3.1 million SNPs. *Nature* **449**: 851–U3.
- 549 FRIEL, N., and H. RUE, 2007 Recursive computing and simulation-free inference for general  
550 factorizable models. *Biometrika* **94**: 661–672.
- 551 GARDNER, K. M., P. BROWN, T. F. COOKE, S. CANN, F. COSTA, *et al.*, 2014 Fast and  
552 cost-effective genetic mapping in apple using next-generation sequencing. *G3 (Bethesda)* **4**:  
553 1681–1687.

554 HEFFELFINGER, C., C. A. FRAGOSO, M. A. MORENO, J. D. OVERTON, J. P. MOTTINGER,  
555 *et al.*, 2014 Flexible and scalable genotyping-by-sequencing strategies for population studies.  
556 *BMC Genomics* **15**: 979.

557 HICKEY, J. M., G. GORJANC, R. K. VARSHNEY, and C. NETTELBLAD, 2015 Imputation of  
558 single nucleotide polymorphism genotypes in biparental, backcross, and topcross populations  
559 with a hidden markov model. *Crop Science* **55**: 1934–1946.

560 HOWIE, B. N., P. DONNELLY, and J. MARCHINI, 2009 A flexible and accurate genotype im-  
561 putation method for the next generation of genome-wide association studies. *PLoS Genetics*  
562 **5**: e1000529.

563 HUANG, B. E., C. RAGHAVAN, R. MAULEON, K. W. BROMAN, and H. LEUNG, 2014 Ef-  
564 ficient imputation of missing markers in low-coverage genotyping-by-sequencing data from  
565 multiparental crosses. *Genetics* **197**: 401–404.

566 KIM, C., H. GUO, W. KONG, R. CHANDNANI, L.-S. SHUANG, *et al.*, 2016 Application of  
567 genotyping by sequencing technology to a variety of crop breeding programs. *Plant Science*  
568 **242**: 14–22.

569 KOVER, P. X., W. VALDAR, J. TRAKALO, N. SCARCELLI, I. M. EHRENREICH, *et al.*, 2009  
570 A multiparent advanced generation inter-cross to fine-map quantitative traits in *Arabidopsis*  
571 *thaliana*. *PLoS Genetics* **5**: e1000551.

572 LI, W., and J. FREUDENBERG, 2009 Two-parameter characterization of chromosome-scale  
573 recombination rate. *Genome Research* **19**: 2300–2307.

574 LI, Y., C. J. WILLER, J. DING, P. SCHEET, and G. R. ABECASIS, 2010 MaCH: Using se-  
575 quence and genotype data to estimate haplotypes and unobserved genotypes. *Genetic Epi-*  
576 *demiology* **34**: 816–834.

577 MACKAY, I. J., P. BANSEPT-BASLER, T. BARBER, A. R. BENTLEY, J. COCKRAM, *et al.*,  
578 2014 An eight-parent multiparent advanced generation inter-cross population for winter-sown  
579 wheat: creation, properties, and validation. *G3 (Bethesda)* **4**: 1603–1610.



580 MARCHINI, J., and B. HOWIE, 2010 Genotype imputation for genome-wide association stud-  
581 ies. *Nature Reviews. Genetics* **11**: 499–511.

582 MONEY, D., K. GARDNER, Z. MIGICOVSKY, H. SCHWANINGER, G. Y. ZHONG, *et al.*, 2015  
583 LinkImpute: fast and accurate genotype imputation for nonmodel organisms. *G3 (Bethesda)*  
584 **5**: 2383–2390.

585 PEI, Y. F., J. LI, L. ZHANG, C. J. PAPASIAN, and H. W. DENG, 2008 Analyses and compari-  
586 son of accuracy of different genotype imputation methods. *PLOS ONE* **3**: e3551.

587 RABINER, L., 1989 A tutorial on hidden markov models and selected applications in speech  
588 recognition. *Proceedings of the IEEE* **77**: 257–286.

589 ROSHYARA, N. R., K. HORN, H. KIRSTEN, P. AHNERT, and M. SCHOLZ, 2016 Compar-  
590 ing performance of modern genotype imputation methods in different ethnicities. *Scientific*  
591 *Reports* **6**: 34386.

592 RUTKOSKI, J. E., J. POLAND, J. L. JANNINK, and M. E. SORRELLS, 2013 Imputation of  
593 unordered markers and the impact on genomic selection accuracy. *G3 (Bethesda)* **3**: 427–  
594 439.

595 SANNEMANN, W., B. E. HUANG, B. MATHEW, and J. LEON, 2015 Multi-parent advanced  
596 generation inter-cross in barley: high-resolution quantitative trait locus mapping for flowering  
597 time as a proof of concept. *Molecular Breeding* **35**: 86.

598 SPINDEL, J., M. WRIGHT, C. CHEN, J. COBB, J. GAGE, *et al.*, 2013 Bridging the genotyp-  
599 ing gap: using genotyping by sequencing (GBS) to add high-density SNP markers and new  
600 value to traditional bi-parental mapping and breeding populations. *Theoretical and Applied*  
601 *Genetics* **126**: 2699–2716.

602 SWARTS, K., H. H. LI, J. A. R. NAVARRO, D. AN, M. C. ROMAY, *et al.*, 2014 Novel methods  
603 to optimize genotypic imputation for low-coverage, next-generation sequence data in crop  
604 plants. *Plant Genome* **7**: 1–12.

605 THEPOT, S., G. RESTOUX, I. GOLDRINGER, F. HOSPITAL, D. GOUACHE, *et al.*, 2015 Effi-  
606 ciently tracking selection in a multiparental population: the case of earliness in wheat. *Ge-*  
607 *netics* **199**: 609–621.

608 WHALEN, A., G. GORJANC, R. ROS-FREIXEDES, and J. M. HICKEY, 2017 Assessment of the  
609 performance of different hidden markov models for imputation in animal breeding. *bioRxiv*  
610 : 227157.

611 WOLFRAM RESEARCH, I., 2016 *Mathematica*. Wolfram Research, Inc., Champaign, Illinois,  
612 version 11.0 edition.

613 XIE, W. B., Q. FENG, H. H. YU, X. H. HUANG, Q. A. ZHAO, *et al.*, 2010 Parent-independent  
614 genotyping for constructing an ultrahigh-density linkage map based on population sequenc-  
615 ing. *Proceedings of the National academy of Sciences of the United States of America* **107**:  
616 10578–10583.

617 ZHENG, C., 2015 Modeling X-linked-linked ancestral origins in multiparental populations. *G3*  
618 (Bethesda) **5**: 777–801.

619 ZHENG, C., M. P. BOER, and F. A. VAN EEUWIJK, 2014 A general modeling framework for  
620 genome ancestral origins in multiparental populations. *Genetics* **198**: 87–101.

621 ZHENG, C., M. P. BOER, and F. A. VAN EEUWIJK, 2015 Reconstruction of genome ancestry  
622 blocks in multiparental populations. *Genetics* **200**: 1073–1087.

623 ZHENG, C., M. P. BOER, and F. A. VAN EEUWIJK, 2018 Recursive algorithms for modeling  
624 genome blocks in a fixed pedigree. Submitted to *G3 (Bethesda)* .

625 ZHENG, C., R. E. VOORRIPS, J. JANSEN, C. A. HACKETT, J. HO, *et al.*, 2016 Probabilistic  
626 multilocus haplotype reconstruction in outcrossing tetraploids. *Genetics* **203**: 119–131.

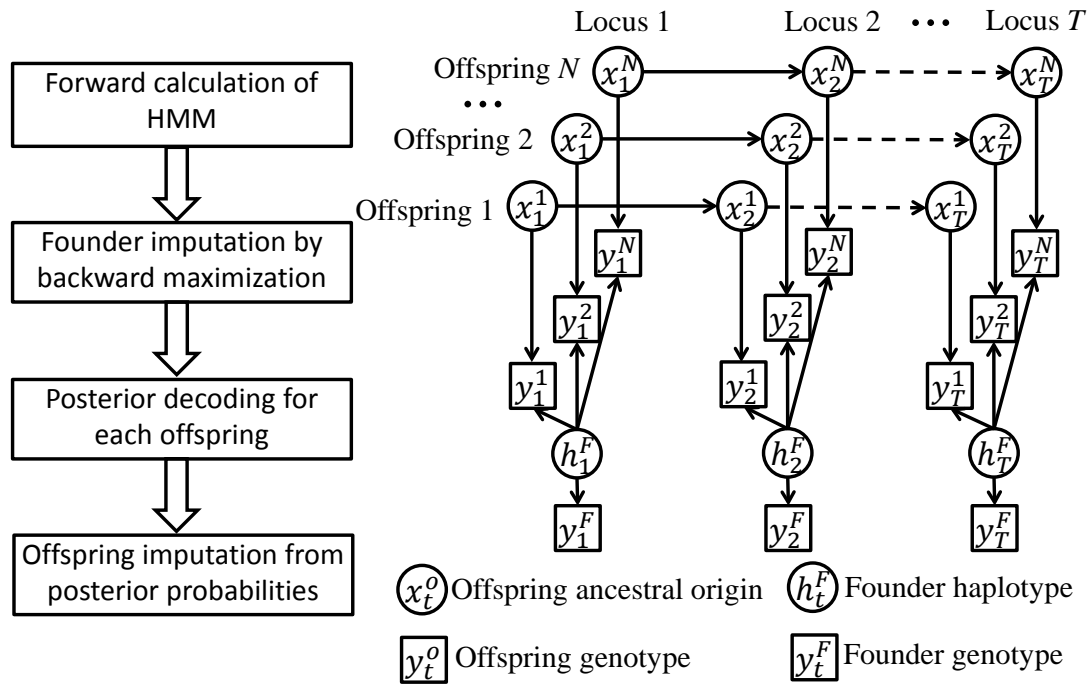
**Figure 1**

Figure 1: Overview of the imputation algorithm. The left panel shows the diagram of magicImpute. The right panel shows the directed acyclic graph of the HMM for  $N$  offspring at  $T$  loci, where the arrows denote probabilistic relationships that are described in the method section. See Table 1 for the symbols in the right panel. In the left panel, the second step of founder imputation results in the estimate of  $h_t^F$  and the third step of posterior decoding results in the posterior probability of  $x_t^o$ , conditional on genotypic data  $y_t^o$  and  $y_t^F$  for  $t = 1 \dots T$  and  $o = 1 \dots N$ .

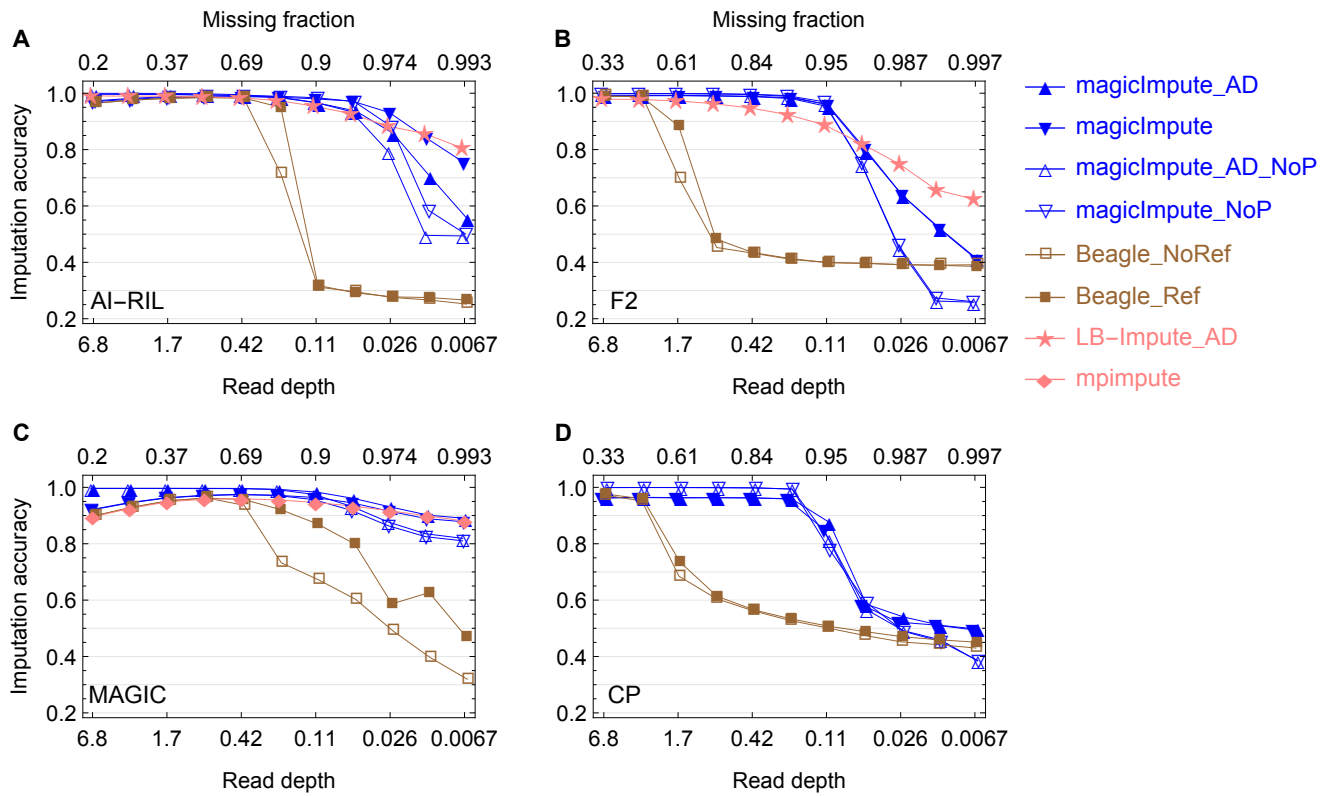


Figure 2

Figure 2: Simulation evaluation on the accuracy of imputing offspring genotypes. Panels A-D show the results for the AI-RIL, the F2, the MAGIC, and the CP, respectively. In the figure legend on the right side, "\_AD" denotes that the input data are allelic depths rather than called genotypes, "\_NoP" denotes that the first two founders' genotypes are not available, and "\_Ref" and "\_NoRef" denotes whether Beagle uses founder haplotypes as reference panels or not. When the input data are called genotypes, complete homozygosity is assumed for the AI-RIL and the MAGIC, and thus their missing fractions on the top axes are smaller than those of the F2 and the CP at the same depths.

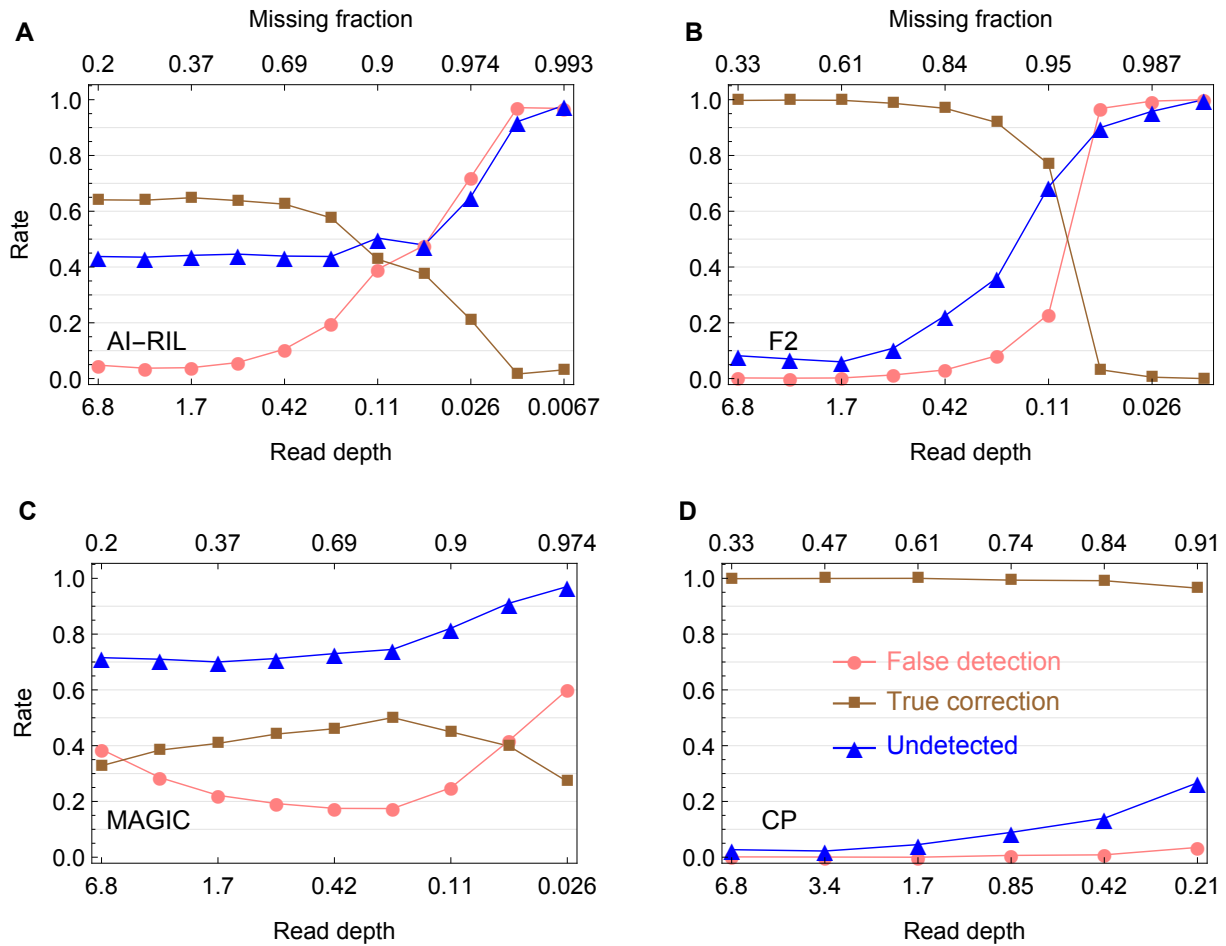


Figure 3

Figure 3: Simulation evaluation on the error detection in offspring genotypes. Panels A-D show the results for the AI-RIL, the F2, the MAGIC, and the CP, respectively, which are obtained by magicImpute with the first two founders' genotypes being unavailable and the input data being called genotypes. The false detection rate (●) denotes the percentage of estimated suspicious genotype errors being not true errors, the true correction rate (■) denotes the percentage of estimated suspicious genotype errors being true and being corrected into the true genotypes, and the undetected rate (▲) denotes the percentage of true genotype errors being not detected.

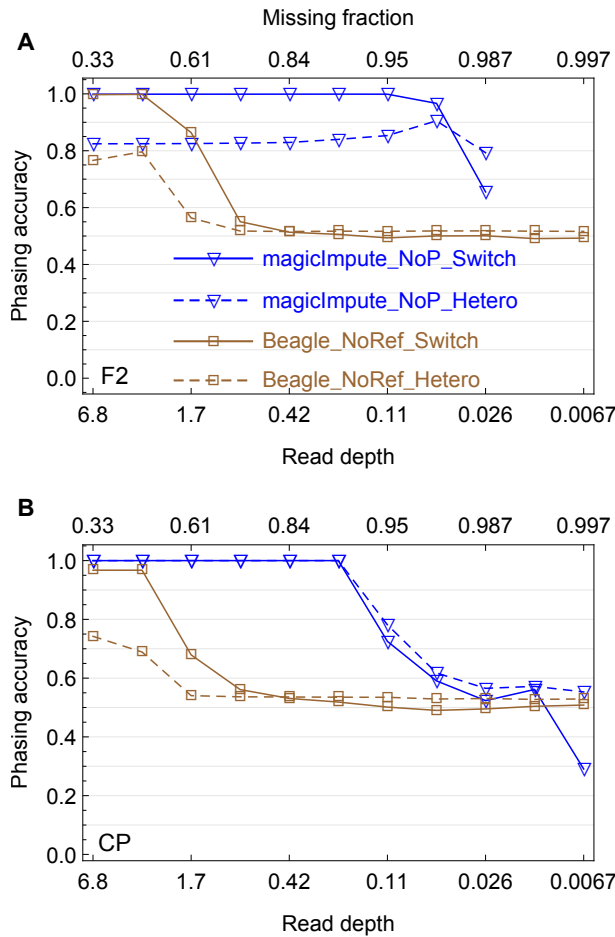


Figure 4

Figure 4: Simulation evaluation on the offspring genotype phasing. Panels A and B show the results obtained by magicImpute and Beagle for the F2 and the CP, respectively. For magicImpute, the first two founders' genotypes are unavailable ("\_NoP"), and for Beagle there are no reference panels ("\_NoRef"). The solid lines denote the switch accuracy ("\_Switch"), one minus the percentage of switch errors to obtain the true haplotype phase; the dashed lines denote the percentage of correctly phased heterozygous genotypes ("\_Hetero").

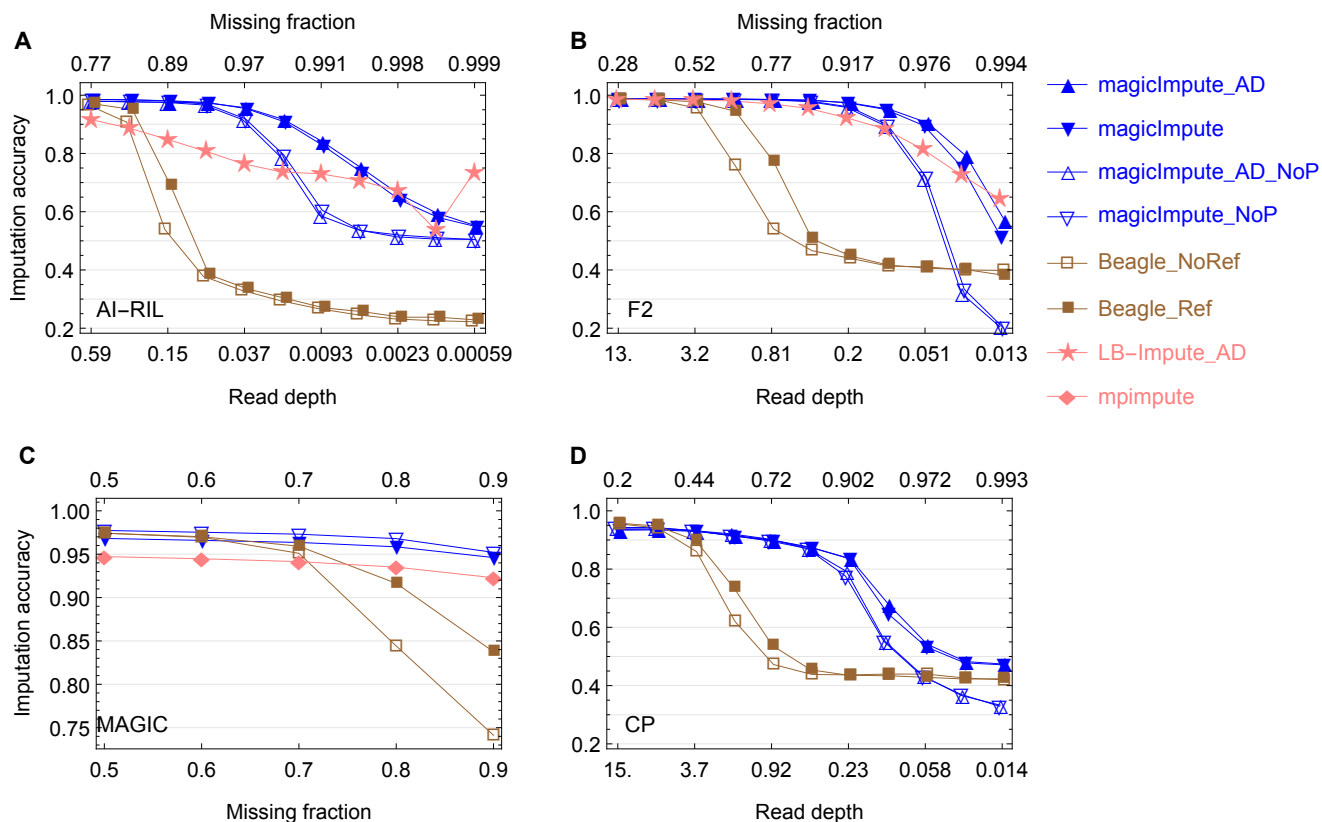


Figure 5

Figure 5: The accuracy of imputing offspring genotypes from real data. Panels A-D show the results for the AI-RIL, the F2, the MAGIC, and the CP, respectively. The figure legend on the right side is the same as that of Figure 2. Allelic depth data are not available for the MAGIC. The extreme large missing fraction or low read depth shows how genotype imputation approaches random imputation with decreasing amount of the input data. In panel A, the large variation of imputation accuracy of LB-Impute at low read depths is due to the corresponding imputation fraction being close to 0 (Figure S8B).

Table 1: List of symbols and their brief descriptions

Symbol	Description
$n_F$	Number of founders
$N$	Number of offspring
$T$	Number of markers (loci)
$h_t^F$	Hidden founder haplotype at locus $t$
$\mathbf{h}^F$	Hidden founder haplotype matrix $\mathbf{h}^F = \{h_t^F\}_{t=1\dots T}$
$x_{ti}$	Hidden ancestral origins at locus $t$ in offspring $i$
$x_{ti}^m, x_{ti}^p$	$x_{ti} = (x_{ti}^m, x_{ti}^p)$ on maternally ( $m$ ) or paternally ( $p$ ) derived chromosome
$d_{ti}$	Genotype at locus $t$ in offspring $i$ that is completely determined by $x_{ti}$ and $h_t^F$
$z_{ti}$	Hidden true genotype at locus $t$ in offspring $i$
$y_{ti}$	Observed genotype at locus $t$ in offspring $i$
$\mathbf{y}^O$	Observed offspring genotype matrix $\mathbf{y}^O = \{y_{ti}\}_{t=1\dots T, i=1\dots N}$
$\mathbf{y}_t^F$	Observed genotypes for all founders at locus $t$
$\mathbf{y}^F$	Observed founder genotype matrix $\mathbf{y}^F = \{y_t^F\}_{t=1\dots T}$
$1U, 2U, UU$	Genotypes containing uncertain allele $U$
$r_1, r_2$	Number of reads for alleles 1 or 2
$\epsilon_O$	Allelic error probability for offspring, independent of read depths
$\epsilon_F$	Allelic error probability for founders, independent of read depths
$phred$	Phred quality score
$\epsilon$	Sequencing error probability $\epsilon = 10^{-phred/10}$
$\pi(x_{1i})$	Prior probability of $x_{1i}$ at locus 1 in offspring $i$
$P(x_{ti} x_{t-1,i})$	Prior transition probability from $x_{t-1,i}$ to $x_{ti}$
$l_{ti}$	$l_{ti} = P(y_{ti} h_t^F, x_{ti}, \epsilon_O, \epsilon_F, \epsilon)$ likelihood at locus $t$ in offspring $i$
$\alpha(x_{ti} h_t^F)$	Posterior probability of $x_{ti}$ conditional on $h_t^F$ and genotypic data from loci 1 to $t$
$\tilde{\alpha}(x_{ti} h_t^F)$	Unnormalized conditional posterior probability of $x_{ti}$
$\gamma(h_t^F)$	Posterior probability of $h_t^F$ conditional on genotypic data from loci 1 to $t$
$\hat{h}_t^F, \hat{x}_{ti}, \hat{z}_{ti}$	Hats denote maximum likelihood estimates
$P_{call}$	Single genotype call if probability of most probable genotype $>$ threshold $P_{call}$
$P_{impute}$	Impute if probability of most probable genotype is $>$ threshold $P_{impute}$
$P_{detect}$	Correct if probability of most probable genotype $>$ threshold $P_{detect}$



Table 2: The running time (in seconds) of genotype imputation for the four real datasets.

Population	Maize AI-RIL	Maize F2	Rice MAGIC	Apple CP
Number of SNPs	13,912	127,059	37,240	13,493
Founder type	inbred	inbred	inbred	outbred
Offspring type	inbred	outbred	inbred	outbred
Number of founders	2	2	8	2
Number of offspring	275	87	178	87
magicImpute	784	212	3170	627
Beagle v4.1	178	31	445	39
LB-Impute	3698	3579	NA	NA
mpimpute	NA	NA	406	NA

## File S1

628

### 629 **Running setups of imputation packages**

#### 630 **magicImpute**

631 The Mathematica command line of magicImpute is given by

```
632 magicImpute[inputfile, model, popdesign, options]
```

633 where `inputfile` specifies the input genotypic data. Here `model` is set to be {"depModel",  
634 "jointModel"} for the population types AI-RIL and MAGIC, so that "depModel" is used  
635 for parental imputation, and "jointModel" is used for offspring imputation. And it is set to  
636 be "jointModel" for the population types F2 and CP, so that "jointModel" is used for  
637 both parental imputation and offspring imputation. See the online manual for details.

638 `popdesign` specifies the breeding design information that is used to compute the process  
639 parameter values of the HMM. For the F2, it is set to be {"Pairing", "Selfing"}. For the  
640 AI-RIL, it is set to be {"RM1-NE-1000", "RM1-NE", ..., "RM1-NE", "Selfing",  
641 ..., "Selfing"} where "RM1-NE" is repeated for 5 times, and "Selfing" is repeated  
642 for 6 times. For the MAGIC, it is set to be {"Pairing", "Pairing", "Pairing",  
643 "Selfing", ..., "Selfing"} where "Selfing" is repeated for 4 times. For the CP,  
644 it is specified in terms of a pedigree file.

645 There are many options for magicImpute. The option `imputingTarget -> All` so  
646 that we by default impute both founder and offspring. The options `founderAllelicError`  
647 `-> 0.005` and `offspringAllelicError -> 0.005` specify  $\varepsilon_F$  and  $\varepsilon_O$ , respectively.  
648 The option `isFounderInbred -> True` specifies that the founders are inbred for the F2,  
649 the AI-RIL, and the MAGIC, and `isFounderInbred -> False` is used for the CP. The  
650 option `imputingThreshold -> 0.9` specifies  $P_{impute}$ . The option `detectingThreshold`  
651 `-> 0.9` specifies  $P_{detect}$ . The option `minPhredQualScore -> 30` specifies that the qual-  
652 ity score  $phred$  so that  $\varepsilon = 10^{-phred/10}$ . The option `priorFounderCallThreshold ->`  
653 `0.99` specifies the prior genotype calling threshold  $P_{call}$  when the input parental data are allelic  
654 depths.

## 655 **Beagle v4.1**

656 The command line used for Beagle v4.1 is given by

```
657 java -jar beagle.21Jan17.6cc.jar ne=100
```

658 where the effective population size is fixed to be 100. In addition, The *gt* option is used to  
659 specify input offspring genotype data, and the *ref* option is used to specify the imputed phased  
660 founder genotypes as the reference panel. We run Beagle with and without the reference panel.

## 661 **LB-Impute**

662 The command line used for LB-Impute is given by

```
663 java -jar LB-Impute.jar -method impute -readerr 0.001  
664     -genotypeerr 0.01 -recombdist 10000000 -window 7  
665     -parentimpute -offspringimpute
```

666 Here the *-readerr* option specifies the sequencing error, and it is set to be 0.001 corresponding  
667 to the quality score 30. The *-genotypeerr* option specifies the genotype error to be 0.01,  
668 corresponding to the depth-independence allelic error probability of 0.005 in magicImpute.  
669 The two founder names are specified by the *-parents* option, and the input and output files are  
670 specified by the options *-f* and *-o*, respectively.

## 671 **mpimpute**

672 The R command line used for mpimpute is given by

```
673 mpimpute(object, what="both", threshold=0.5, calls="discrete")
```

674 Here the *what* option is set so that we impute both founders and offspring, and input genotypic  
675 data and pedigree information are specified by the *object*.

676 **Supplementary figures**

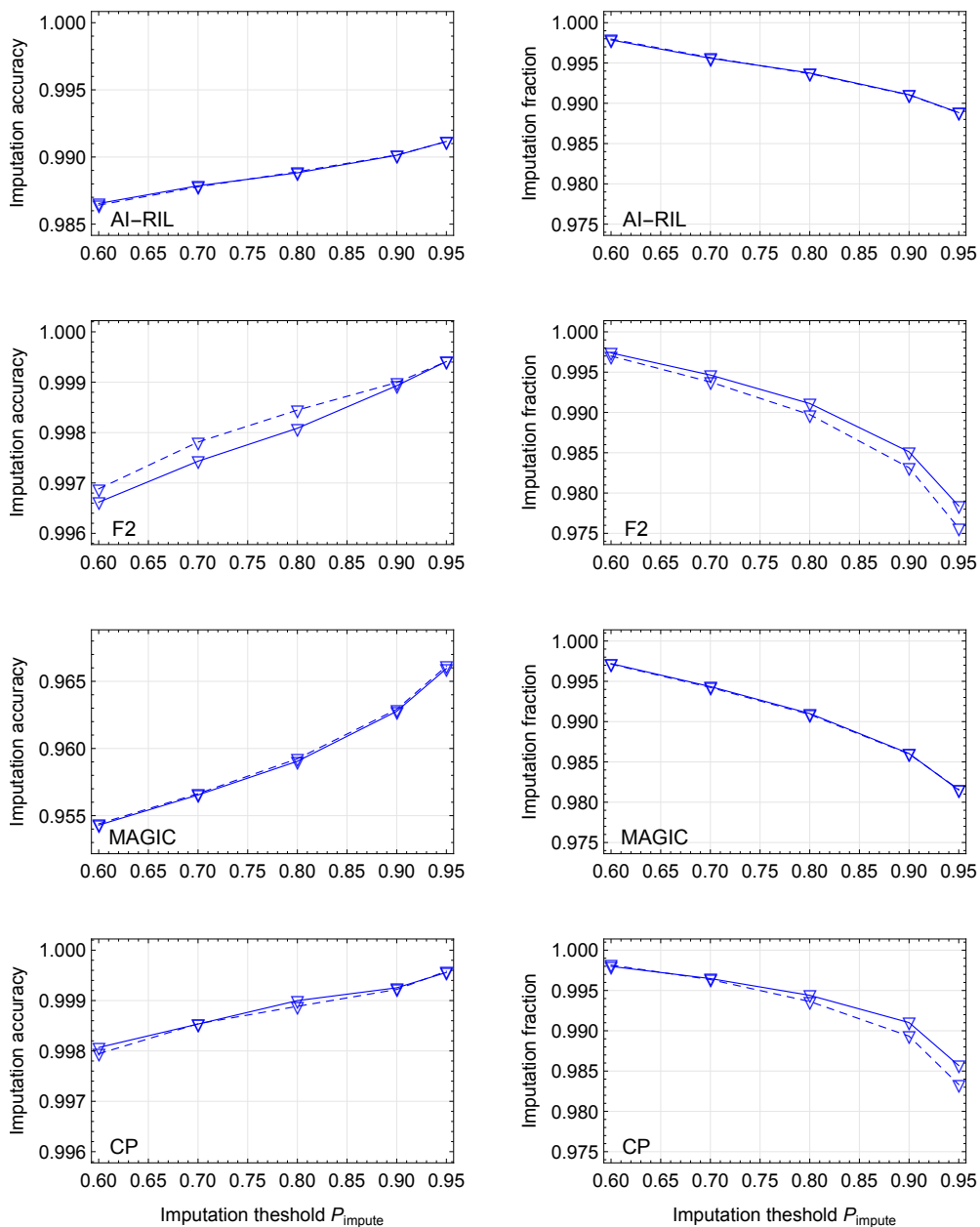


Figure S1: Sensitivity analysis of imputation threshold  $P_{impute}$  for the algorithm `magicImpute`. Panels from top to bottom denote the results for the AI-RIL, the F2, the MAGIC, and the CP, respectively. The solid and dashed lines denote the results corresponding to input parameter  $\epsilon_O = 0.005$  and  $0.05$ , respectively. The left and right panels denote the results for imputation accuracy and imputation fraction, respectively, which are obtained from the simulated datasets with the input data being called genotypes at read depth 0.85 and the first two founders' genotypes being not available.

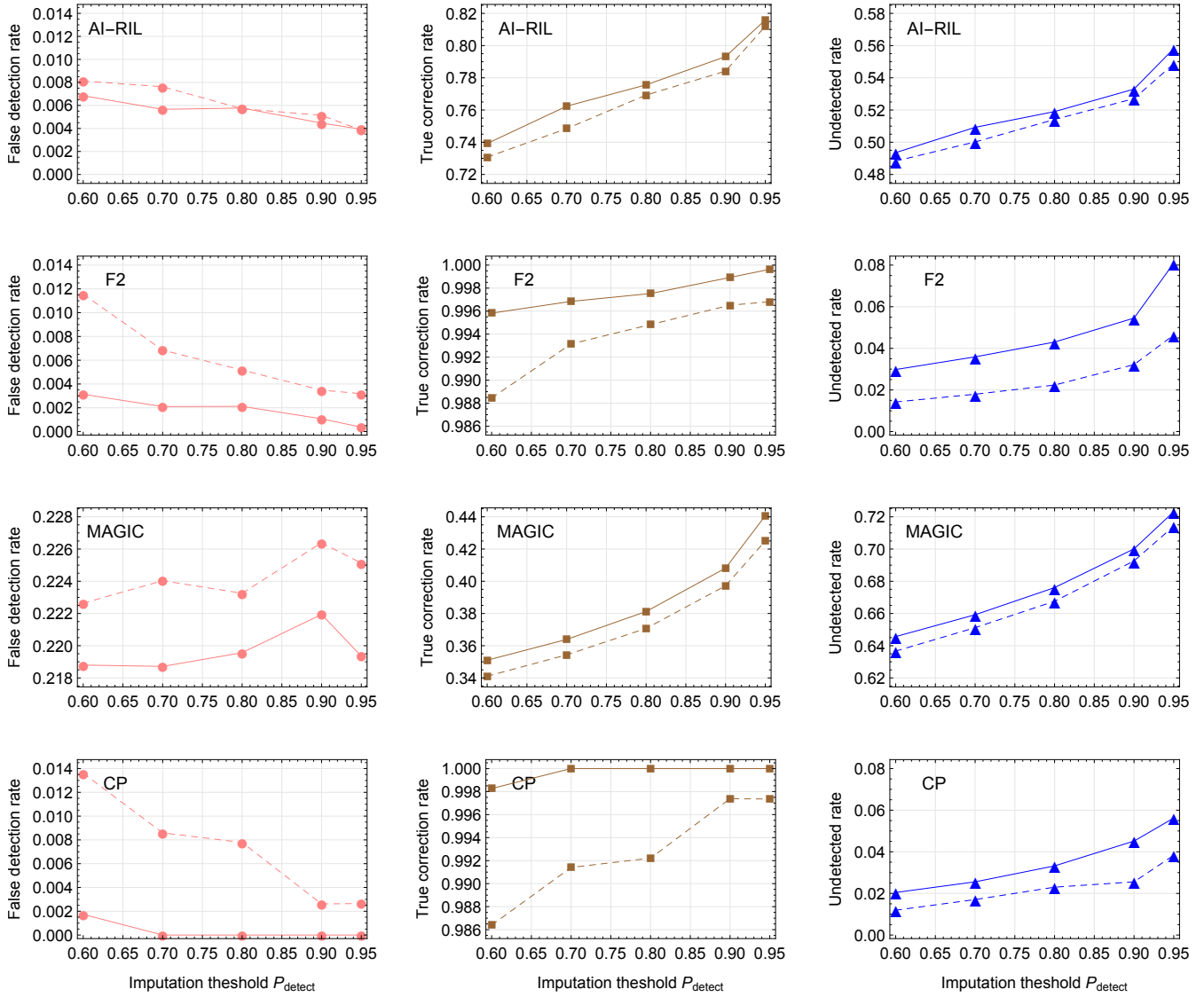


Figure S2: Sensitivity analysis of error detection threshold  $P_{detect}$  for the algorithm magicImpute. Panels from top to bottom denote the results for the AI-RIL, the F2, the MAGIC, and the CP, respectively. The solid and dashed lines denote the results corresponding to input parameter  $\epsilon_O=0.005$  and  $0.05$ , respectively. The left, middle and right panels denote false detection rate, true correction rate, and undetected rate, respectively. For each simulated dataset with population size 200 and read depth 0.85, the results are obtained with the input data being called genotypes and the first two founders' genotypes being not available.

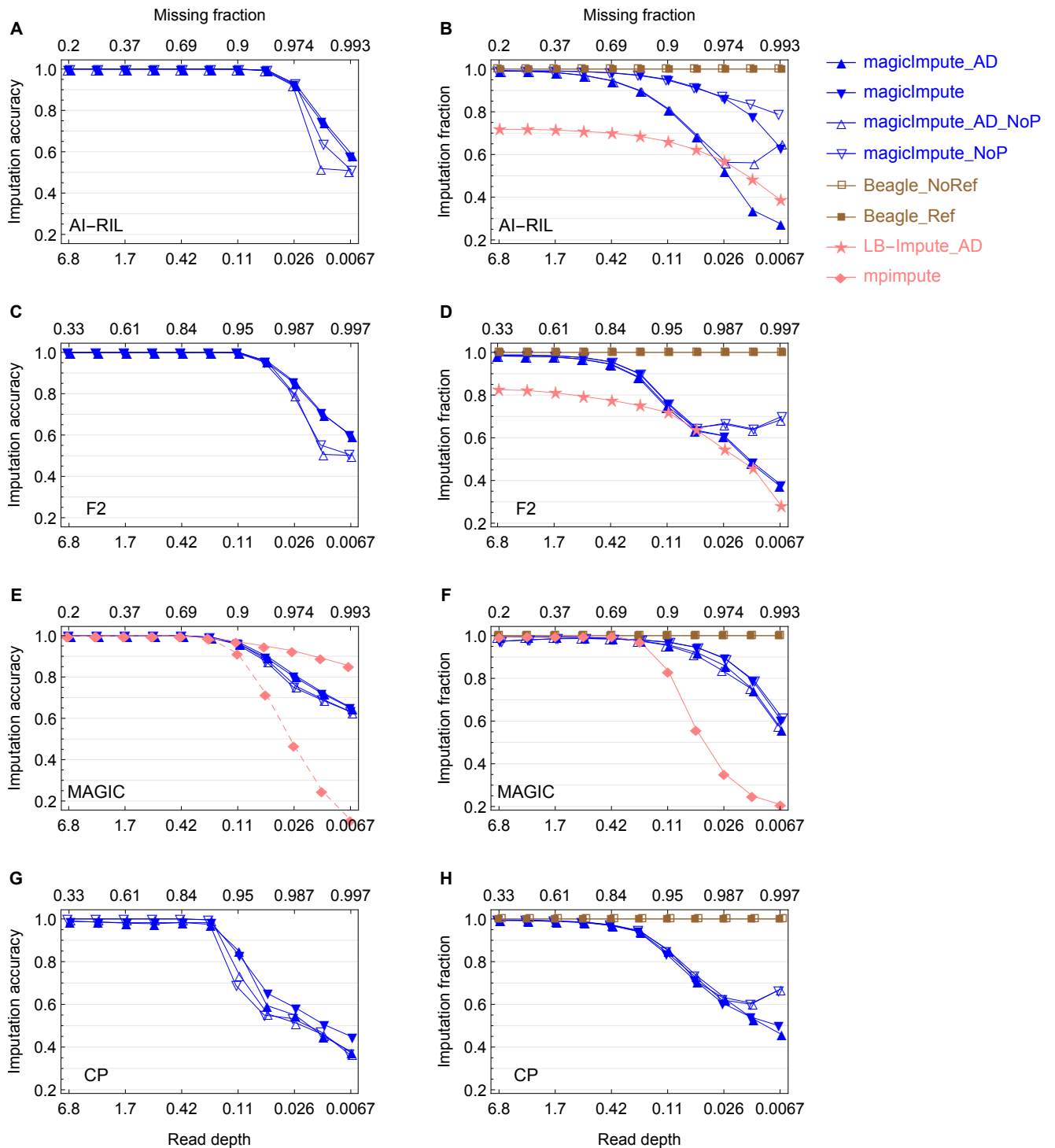


Figure S3: Simulation evaluation on the accuracy of imputing founder genotypes (left panels) and imputation fraction of offspring genotypes (right panels). Panels A&B, C&D, E&F, and G&H denote the results for the AI-RIL, the F2, the MAGIC, and the CP, respectively. The dashed lines in panel E denotes the mpimpute imputation fraction of founder genotypes. Beagle and LB-Impute do not impute founder genotypes, and magicImpute always imputes all the founder genotypes.

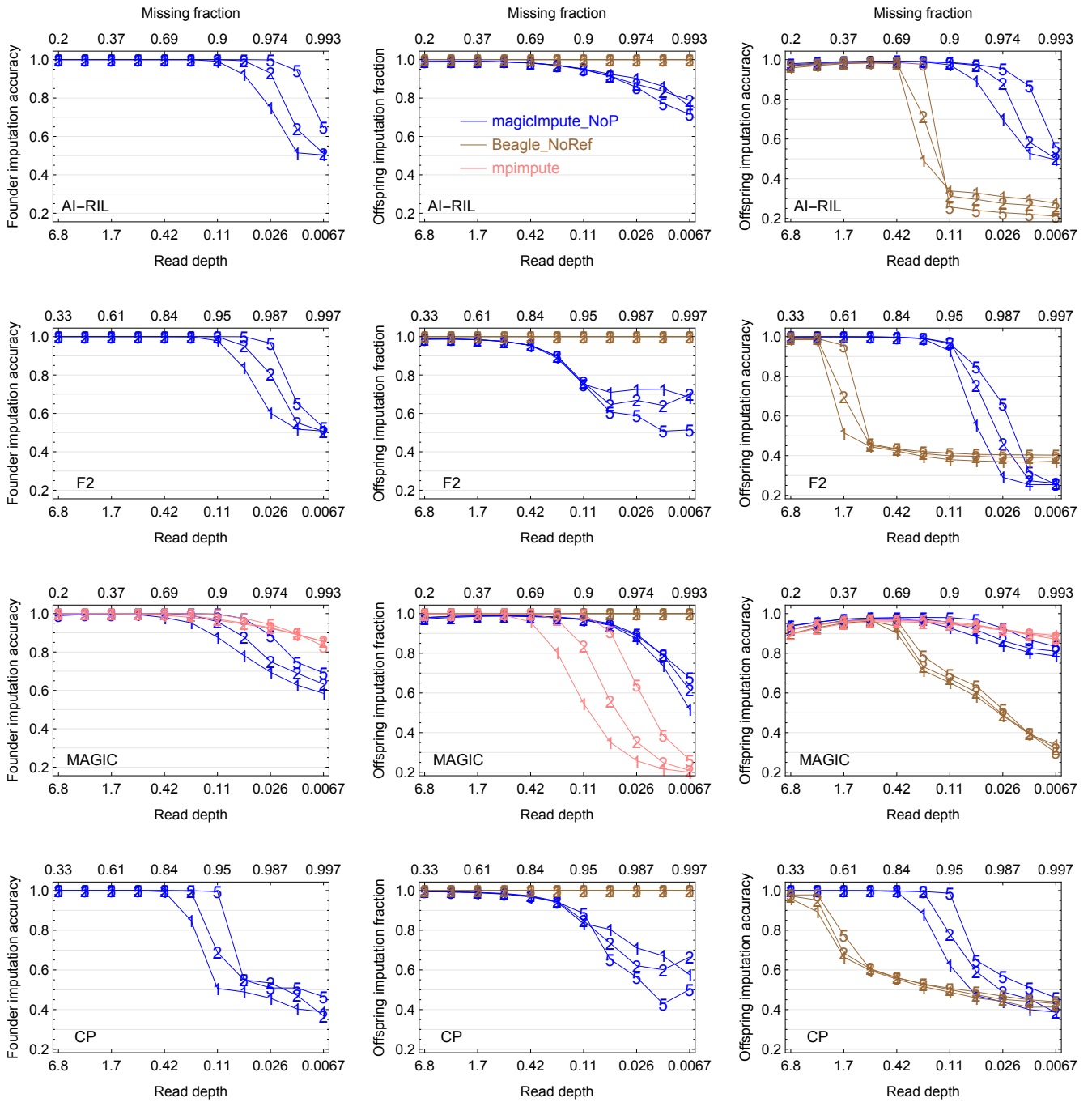


Figure S4: Dependencies of genotype imputation on population size. Panels from top to bottom denote the results for the AI-RIL, the F2, the MAGIC, and the CP, respectively. For magicImpute, the input data are called genotypes and the first two founders are missing; no reference panels for Beagle imputation. The plot markers "1", "2", and "5" denote population sizes 100, 200, and 500, respectively.

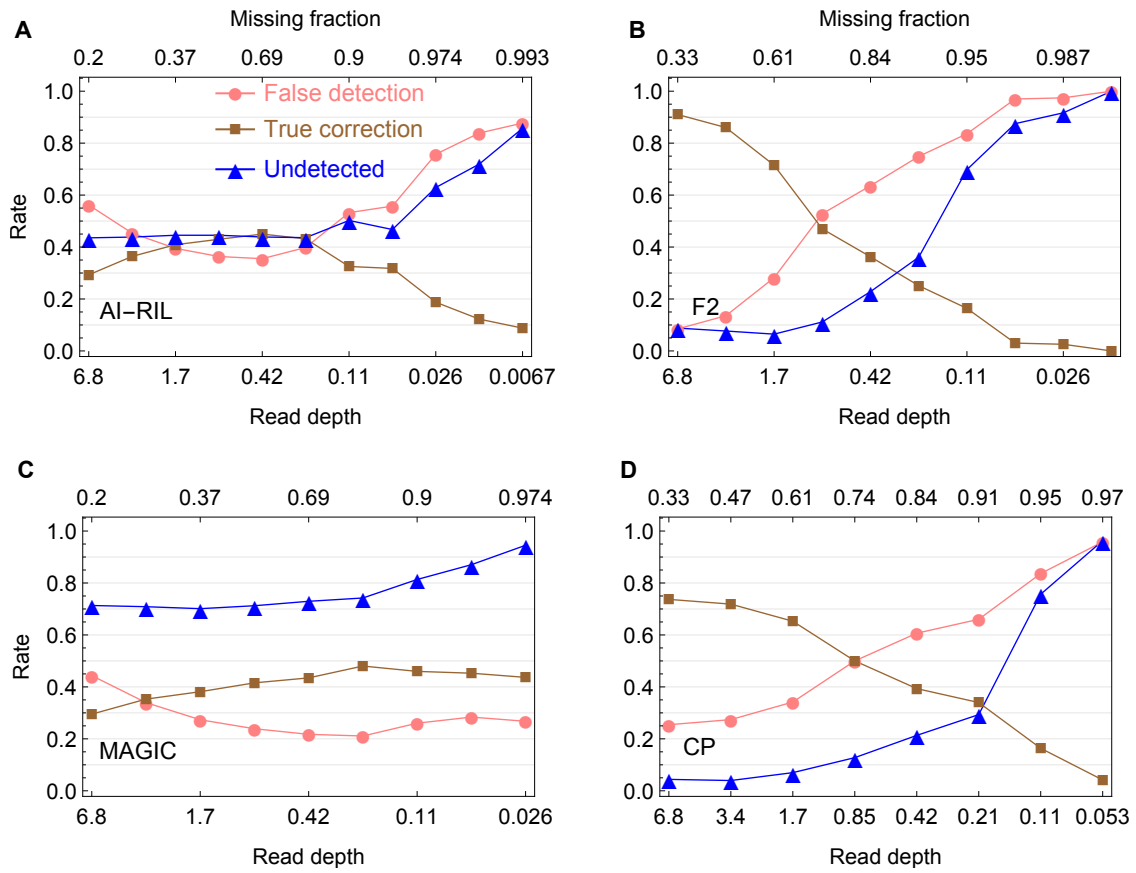


Figure S5: Similar to Figure 3 for the error detection by magicImpute but with the first two founders' genotypes being available.



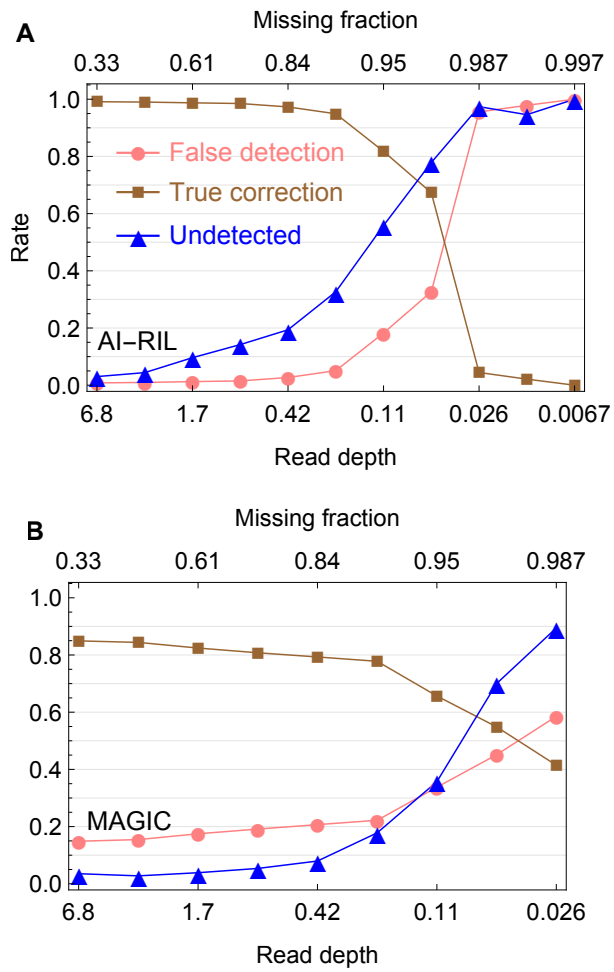


Figure S6: Similar to Figure 3 for the error detection by magicImpute but without assuming homozygosity for the almost homozygous populations AI-RIL (A) and MAGIC (B).

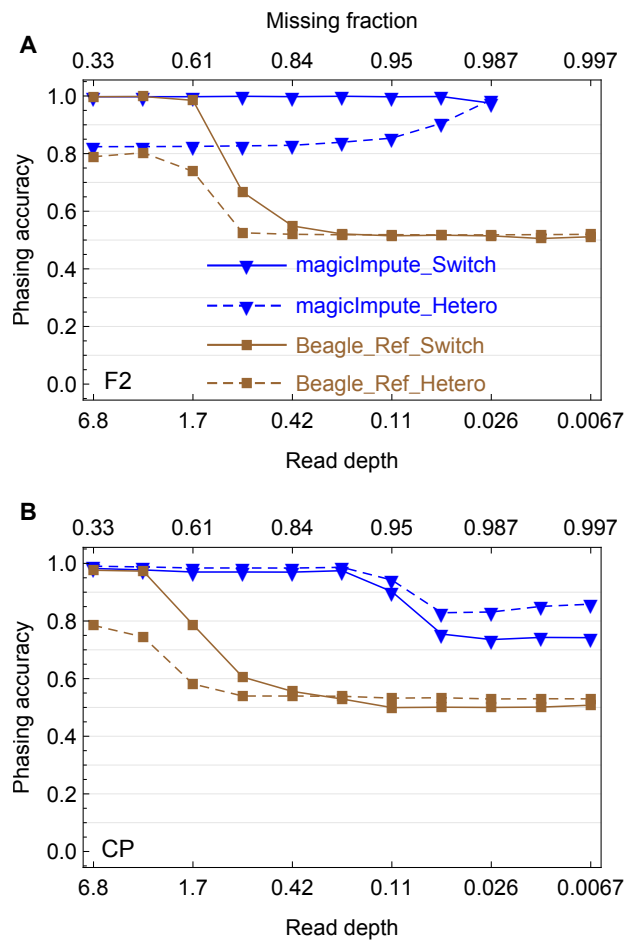


Figure S7: Similar to Figure 4 for the offspring phasing but for magicImpute with the first two founders' genotypes being available and for Beagle with the founder haplotypes (imputed by magicImpute) being the reference panels.

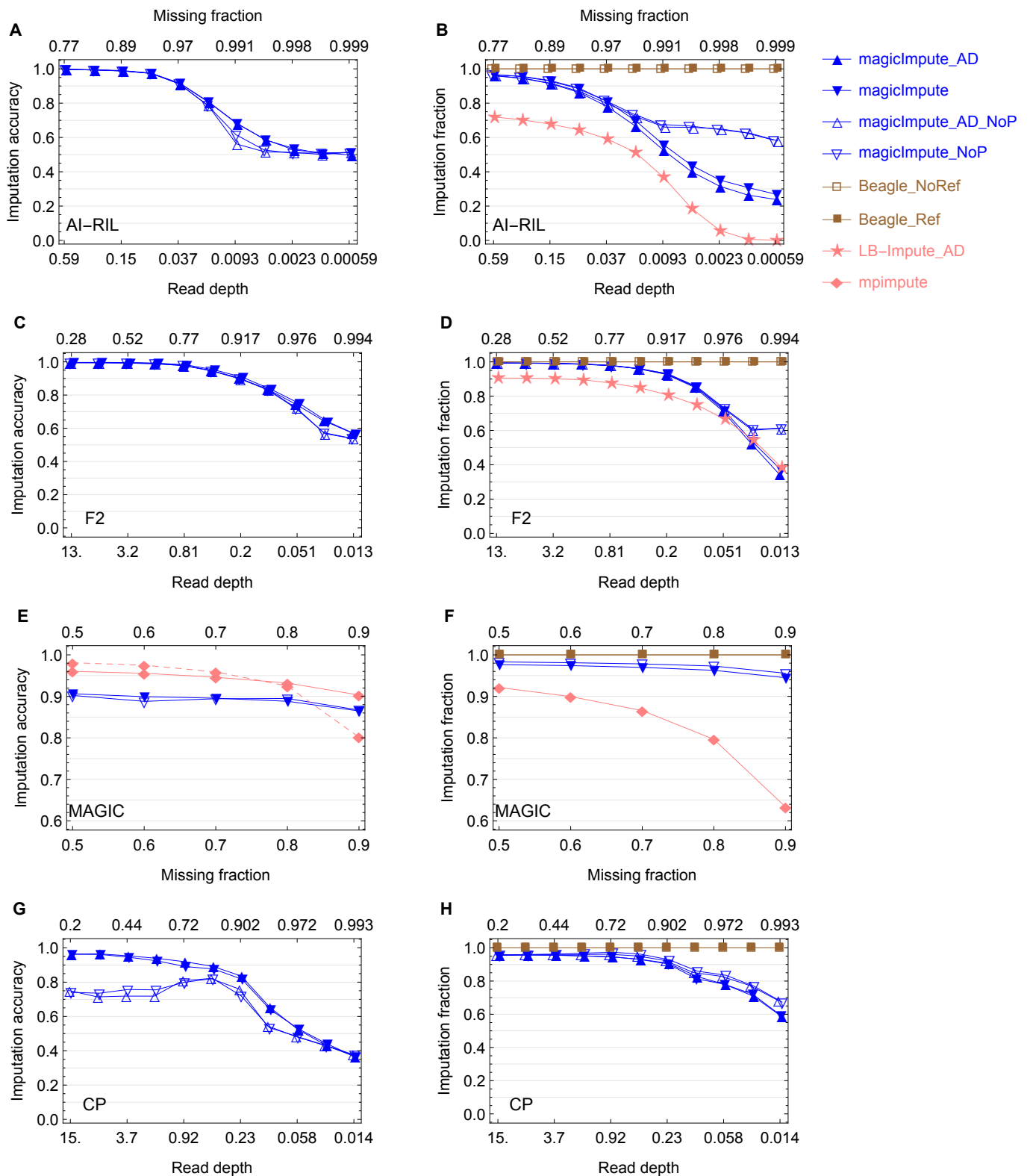


Figure S8: Evaluation on the accuracy of imputing founder genotypes and imputation fraction of offspring genotypes by real data. Panels A&B, C&D, E&F, and G&H denote the results for the AI-RIL, the F2, the MAGIC, and the CP, respectively. The dashed lines in panel E denotes the mpimpute imputation fraction of founder genotypes.