

1 **A heuristic method for fast and accurate phasing and**
2 **imputation of single nucleotide polymorphism data in bi-**
3 **parental plant populations**

4

5 Serap Gonen, Valentin Wimmer, R. Chris Gaynor, Ed Byrne, Gregor Gorjanc, John
6 M. Hickey*

7

8 S. Gonen, G. Gorjanc, R.C. Gaynor and J.M. Hickey The Roslin Institute and Royal
9 (Dick) School of Veterinary Studies, University of Edinburgh, Easter Bush Research
10 Centre, Midlothian EH25 9RG, UK

11 V. Wimmer KWS SAAT SE, Grimsehlstr. 31, 37574 Einbeck, Germany

12 E. Byrne KWS-UK Ltd, 56 Church Street, Thriplow, Hertfordshire, SG8 7RE, UK

13 Received _____ *Corresponding author (john.hickey@roslin.ed.ac.uk)

14

15 **Abbreviations:** LD, low-density; HD, high-density; SNP, single nucleotide
16 polymorphism; cM, centiMorgan.

17 **Abstract**

18 This paper presents a new heuristic method for phasing and imputation of
19 genomic data in diploid plant species. Our method, called AlphaPlantImpute,
20 explicitly leverages features of plant breeding programs to maximise the accuracy of
21 imputation. The features are a small number of parents, which can be inbred and
22 usually have high-density genomic data, and few recombinations separating parents
23 and focal individuals genotyped at low-density (i.e. descendants that are the
24 imputation targets). AlphaPlantImpute works roughly in three steps. First, it identifies
25 informative low-density genotype markers in parents. Second, it tracks the inheritance
26 of parental alleles and haplotypes to focal individuals at informative markers. Finally,
27 it uses this low-density information as anchor points to impute focal individuals to
28 high-density.

29 We tested the imputation accuracy of AlphaPlantImpute in simulated bi-
30 parental populations across different scenarios. We also compared its accuracy to
31 existing software called PlantImpute. In general, AlphaPlantImpute had better or
32 equal imputation accuracy as PlantImpute. The computational time and memory
33 requirements of AlphaPlantImpute were tiny compared to PlantImpute. For example,
34 accuracy of imputation was 0.96 for a scenario where both parents were inbred and
35 genotyped at 25,000 markers per chromosome and a focal F_2 individual was
36 genotyped with 50 markers per chromosome. The maximum memory requirement for
37 this scenario was 0.08 GB and took 37 seconds to complete.

38 **Introduction**

39 This paper presents a new heuristic method for phasing and imputation of
40 single nucleotide polymorphism (SNP) array data in diploid plant species. High-
41 density SNP array data in plant breeding populations is increasingly valuable for
42 genomic selection and for identifying regions of the genome that underlie traits of
43 interest in genome-wide association studies. The accuracy of genomic selection and
44 power of association studies increases with the number of individuals and with the
45 density of SNP markers. However, the cost of genotyping many individuals at high-
46 density is high. This high cost is a barrier to the adoption of genomic selection in
47 plant breeding programs where the number of selection candidates in each cycle can
48 be very large. An effective strategy to overcome this cost barrier is to genotype a
49 proportion of the population at high-density, phase their genotypes, and use this data
50 for imputation of large numbers of individuals genotyped at low-density (Jacobson et
51 al., 2014, 2015; Gorjanc et al., 2017a; b). This strategy has been widely adopted in
52 livestock and human populations, partly because genotype imputation tools that work
53 well in these populations are widely available (Kong et al., 2008; Howie et al., 2009;
54 Druet and Georges, 2010; Li et al., 2010; Sargolzaei et al., 2011; Hickey et al., 2011;
55 Cleveland and Hickey, 2013; Hickey and Kranis, 2013; VanRaden et al., 2015;
56 O’Connell et al., 2016; Loh et al., 2016; Antolín et al., 2017).

57 Bi-parental populations that are widely used in plant breeding have four
58 features that make them ideal for imputation. First, they are derived from only two
59 parents. High-density genotyping of the two parents and low-density genotyping of
60 focal individuals (i.e., descendants that are the imputation targets) is an effective low-
61 cost strategy in these populations. Second, the number of meioses separating parents

62 and focal individuals is small. This means that parental haplotypes remain largely
63 intact in focal individuals, which simplifies imputation. Third, they have well-known
64 crossing structures that could be informative for imputation, although the process of
65 selfing or the creation of doubled haploids can add complications that are not present
66 in human and livestock settings. However, these “complications” can in certain
67 situations empower imputation. Finally, parents that contribute to a bi-parental
68 population are usually inbred. This means that they are homozygous at many loci and
69 the majority of their genome is phased *de facto*.

70 A recent simulation study demonstrated that achieving high imputation
71 accuracies could empower genomic selection in bi-parental populations (Gorjanc et
72 al., 2017a; b). The high imputation accuracies with SNP array data were achieved
73 using the PlantImpute software (Nettelblad et al., 2009; Hickey et al., 2015). The
74 main drawback of PlantImpute is that it has large computational requirements in
75 terms of time and memory. This makes it impractical for routine use in breeding
76 programs. Existing software for imputation in livestock or human populations do not
77 have large computational requirements. However, software for imputation in livestock
78 or human populations are not designed to leverage features of plant breeding
79 programs, and in some cases, cannot work where selfing and bi-sexuality is common.
80 To our knowledge, existing imputation software for plant breeding programs (e.g.,
81 (Swarts et al., 2014)) are not explicitly designed for imputation of SNP array
82 genotypes in bi-parental populations.

83 This paper presents a new heuristic method, called AlphaPlantImpute, for
84 phasing and imputation of SNP array data in diploid plant species. AlphaPlantImpute
85 works roughly in three steps. First, it identifies markers fully or partially informative

86 for parent-of-origin. Second, it tracks the inheritance of parental alleles and
87 haplotypes to focal individuals at informative markers. Finally, it uses this low-
88 density information as anchor points to impute focal individuals to high-density.

89 We tested the accuracy of AlphaPlantImpute in simulated bi-parental
90 populations across different scenarios. These scenarios varied in the levels of
91 inbreeding in the parents, the number of selfing events separating parents and focal
92 individuals, the chromosome size (i.e. recombination rate) and the number of markers
93 on the low-density array. We calculated the accuracy of imputation within each
94 scenario as the correlation between the true and imputed genotypes. In general,
95 AlphaPlantImpute gave excellent accuracy of imputation and typically outperformed
96 or performed equally as well as PlantImpute for the accuracy of imputation. The
97 computational time and memory requirements of AlphaPlantImpute were always tiny
98 compared to that of PlantImpute.

99 **Materials and methods**

100 *Definitions*

101 A focal individual is an individual that is to be imputed. A fully informative
102 marker is one where the two parents have opposing homozygous genotypes, i.e.,
103 genotypes 0 and 2 (note that the method is agnostic of which allele is the reference
104 allele). A partially informative marker is where one parent is homozygous and the
105 other is heterozygous. Markers where parents are fixed for the same allele or where
106 both parents are heterozygous are uninformative. The high-density (**HD**) array is the
107 array at which parents have genotypes and is the target array for imputation. In our
108 test datasets, the HD array consisted of 25,000 SNP markers. The low-density (**LD**)
109 array is the array at which focal individuals have genotypes. We tested eight LD
110 arrays (see below), all of which were nested subsets of the HD array.

111 *Description of the method*

112 We present a new heuristic method, called AlphaPlantImpute, for phasing and
113 imputation of SNP array data in diploid plant species. In detail, our method has five
114 steps: (1) Identify markers that are informative for parent-of-origin of alleles in focal
115 individuals; (2) Infer the most likely linked alleles at two markers; (3) Phase and
116 assign parent-of-origin for focal individual's alleles; (4) Impute focal individual to
117 high-density using low-density anchors captured in step 3; and (5) Impute markers in
118 recombined regions. Impute markers adjacent to recombination locations. Step 1 is
119 the only step applied to groups of focal individuals together. Steps 2, 3, 4 and 5 are
120 applied for each focal individual separately. A description of the definitions used and

121 of each step is given below and a schematic is given in Figure 1 (a more detailed
122 schematic is given in Supplementary Figure 1).

123 *Method steps*

124 Step 1: Identify informative low-density markers in parents

125 In the first step we determine which low-density markers are fully or partially
126 informative in parents, which is used in the following steps to infer parent-of-origin of
127 phased alleles in focal individuals. For example, in Figure 1 eight of the ten markers
128 on the HD array genotyped in the parents are fully informative and two (markers 2
129 and 9) are uninformative. Of the ten HD markers, five (markers 1, 3, 5, 7, 9) are also
130 on the LD array, which was used to genotype focal individuals. Of these five LD
131 markers, four are informative and one (marker 9) is uninformative.

132 Step 2: Infer the most likely linked alleles at two markers

133 In the second step we infer the most likely linked alleles at two markers for all
134 pairs of informative markers, which is used in the following steps to phase
135 heterozygous markers in focal individuals. If parent haplotypes are inherited directly
136 without recombination, the most likely linked alleles at two markers recover the
137 parent haplotypes. When this is not the case, the most likely linked alleles at two
138 markers indicate a potential recombination hotspot or marker map error for the
139 population. For each pair of informative markers we perform three steps.

140 2a) First, identify focal individuals that are homozygous at the first and the
141 second marker.

142 2b) Second, count the number of times focal individuals have genotype:

- 143 • 0 for the first and 0 for the second marker (diplotype 0-0),
- 144 • 0 for the first and 2 for the second marker (diplotype 0-2),
- 145 • 2 for the first and 0 for the second marker (diplotype 2-0), and
- 146 • 2 for the first and 2 for the second marker (diplotype 2-2).

147 2c) Third, compare the count of 0-0 to 0-2 and of 2-2 to 2-0. If the count of 0-
148 0 is higher than 0-2 and 2-2 is higher than 2-0, then the 0 (1) allele at the first marker
149 is commonly linked to the 0 (1) allele at the second marker. If the count of 0-2 is
150 higher than 0-0 and 2-0 is higher than 2-2, then the 0 (1) allele at the first marker is
151 commonly linked to the 1 (0) allele at the second marker. For example, in Figure 1 2-
152 2 and 0-0 are the two most frequent diplotypes at markers 1 and 3, which suggests the
153 most likely linked alleles are 1-1 and 0-0.

154 Step 3: Phase and assign parent-of-origin for focal individual's alleles

155 In the third step we phase alleles in focal individuals and assign their parent-
156 of-origin. We perform this first for the homozygous markers and then for the
157 heterozygous markers.

158 *3a) Phase homozygous markers*

159 We phase alleles at homozygous markers as the 0 allele for both haplotypes
160 when the genotype is 0 and as the 1 allele when the genotype is 2. For example, in
161 Figure 1 the focal individual ID_Y has genotype 2 for marker 7 and we phase it as the
162 1 allele for both haplotypes.

163 *3b) Assign parent-of-origin to alleles at homozygous markers*

164 We assign parent-of-origin for phased alleles in the step 3a based on the
165 informative markers in the step 1. For example, in Figure 1 marker 7 is informative.
166 At this marker, the Parent_A has the 0 allele, while the Parent_B has the 1 allele.
167 Focal individual ID_Y has genotype 2, which suggests that both of the 1 alleles were
168 inherited from the Parent_B. Focal individual ID_Y is also homozygous at marker 9,
169 with genotype 0, but this marker is not informative and we cannot assign parent-of-
170 origin to phased alleles.

171 *3c) Phase heterozygous marker*

172 We phase alleles at heterozygous markers iteratively based on the most likely
173 linked alleles in the step 2. Specifically, we perform four steps. We start at the first
174 heterozygous marker . For example, in Figure 1 the first marker for which the focal
175 individual ID_Y is heterozygous is marker 1.

176 3c1) First, phase the first heterozygous marker randomly as the 1 allele for the
177 first haplotype and the 0 allele for the second haplotype.

178 3c2) Second, phase the second heterozygous marker based on the the most
179 likely linked alleles in the step 2. For example, in Figure 1 the second heterozygous
180 marker is marker 3. Information from the most likely linked alleles suggest that the 0
181 (1) allele at marker 1 is linked to the 0 (1) allele at marker 3. Using this information,
182 we phase marker 3 alleles of ID_Y as the 1 allele for the first haplotype and the 0
183 allele for the second haplotype. We continue moving from left-to-right until the last
184 heterozygous marker is phased.

185 3c3) Third, we repeat steps 3c1 and 3c2, but this time starting from the last
186 heterozygous marker and progressing to the first heterozygous marker.

187 3c4) Finally, we derive a consensus between the haplotypes derived from
188 moving left-to-right and right-to-left along the chromosome. If they disagree, set the
189 consensus haplotypes to missing. If only one is filled, set the consensus haplotype to
190 the filled information.

191 *3d) Assign parent-of-origin to alleles at heterozygous marker*

192 We assign parent-of-origin for phased alleles in the step 3c based on the
193 informative markers in the step 1. For example, in Figure 1 focal individual ID_Y is
194 heterozygous at marker 1. At this marker, the 1 allele on ID_Y's first haplotype is
195 inherited from Parent_A and the 0 allele on ID_Y's second haplotype is inherited
196 from Parent_B. If the marker is partially informative, we assign both the parent-of-
197 origin and the haplotype-of-origin (i.e., first or second haplotype of the parent that is
198 heterozygous for that marker).

199 Step 4: Impute focal individual to high-density using anchors from the step 3

200 *4a) Fill uninformative homozygous markers*

201 For uninformative homozygous markers at HD that are not genotyped in the
202 focal individual at LD, we phase and impute the focal individual with the parental
203 information. For example, in Figure 1 both parents have genotype 0 for marker 2, so
204 focal individual ID_Y is imputed as genotype 0.

205 *4b) Assign parent-of-origin to HD marker alleles*

206 For markers on the HD array, assign parent-of-origin to marker alleles based
207 on the parent-of-origin assignment of the two nearest marker alleles on the LD array.
208 For example, in Figure 1 marker 6 is not genotyped on the LD array but the two

209 neighbouring markers 5 and 7 are genotyped on the LD array. We have assigned the
210 second haplotype of focal individual ID_Y to Parent_B for both markers 5 and 7. We
211 therefore also assign marker 6 to Parent_B for the second haplotype. We have
212 assigned the first haplotype of focal individual ID_Y to Parent_A for marker 5 and to
213 Parent_B for marker 7. We conclude that there was a potential recombination around
214 marker 6 at the first haplotype and we do not assign parent-of-origin for this allele.

215 *4c) Phase and impute HD markers using parent-of-origin assignment from*
216 *step 4b*

217 For HD markers with assigned parent-of-origin in step 4b, we phase the allele
218 inherited from that parent for the haplotype of the focal individual. If we have phased
219 both alleles at a marker, we impute the genotype as the sum of the two alleles on the
220 two haplotypes of the focal individual. If parent-of-origin has not been assigned for
221 one or both alleles of the focal individual, we leave the genotype as missing.

222 Step 5. Impute markers in recombined regions

223 We phase and impute missing HD markers in potentially recombined regions
224 in one of two ways. We either (1) impute expected genotype dosage as the average of
225 the alleles of the two parents; or (2) phase and impute using information from a
226 genetic or physical map. For (2), we first identify the two closest neighbouring
227 markers that were informative and phased, second use the distance between these two
228 markers as a weight to phase the missing alleles as the weighted average of the alleles
229 of the two parent haplotypes, and third impute expected genotype dosage as in (1).

230 *Implementation*

231 We have implemented the method in a program called AlphaPlantImpute,
232 which is controlled by a specification file that contains some user specified thresholds
233 and the addresses of input files. The required input data are membership of
234 individuals to the bi-parental populations, HD genotypes for parents, and LD
235 genotypes of focal individuals. The output data are imputed genotypes, phased
236 haplotypes, inferred parent-of-origin for focal individual haplotypes, and information
237 on whether a marker is informative. AlphaPlantImpute implements some data editing
238 checks, which are described in the user manual.

239 *Examples of implementation: Description of datasets*

240 To test the imputation accuracy of AlphaPlantImpute, testing datasets of a
241 subset of the scenarios described in Hickey et. al. 2015 were simulated. This enabled
242 the comparison of AlphaPlantImpute with PlantImpute without re-running
243 PlantImpute with its large computational cost. Although the simulation design is
244 largely a replication of that in Hickey et. al. 2015, a brief description of the general
245 structure and simulation method of the different scenarios tested is given below for
246 completeness.

247 *Simulation of genomic data*

248 Sequence data for 100 base haplotypes for a single chromosome were
249 simulated using the Markovian Coalescent Simulator (Chen et al., 2009) and
250 AlphaSim (Faux et al., 2016). The base haplotypes were 10^8 base pairs in length, with
251 a per site mutation rate of 1.0×10^{-8} and a per site recombination rate that varied across
252 scenarios. The different recombination rates simulated were 0.25×10^{-8} , 0.5×10^{-8} ,
253 1.0×10^{-8} , 1.5×10^{-8} , 2.0×10^{-8} , 3.0×10^{-8} , and 4.0×10^{-8} , resulting in chromosome sizes

254 of 25, 50, 100, 150, 200, 300, and 400 centiMorgans (cM), respectively. The effective
255 population size (N_e) was set at specific points during the simulation to mimic changes
256 in N_e in a crop such as maize (*Zea mays L.*). These set points were: 100 in the base
257 generation, 1000 at 100 generations ago, and 10,000 at 2000 generations ago, with
258 linear changes in between. The resulting whole-chromosome haplotypes had
259 approximately 80,000 segregating sites in total.

260 *Simulation of a pedigree*

261 A pedigree of 11,266 individuals was constructed. The pedigree was initiated
262 from six outbred founders (A, B, C, D, E, F). These six founders were crossed to
263 generate the founder bi-parental populations (AxB, CxD, ExF). These founder bi-
264 parental populations were selfed to F_1 , F_2 , F_4 , F_{10} , or F_{20} , resulting in different levels
265 of inbreeding in the parents. To properly propagate the residual heterozygosity in
266 these parents, they were crossed to generate 100 pairs of F_1 individuals. F_1 individuals
267 were selfed to generate 100 F_2 individuals. F_2 individuals were selfed to generate 100
268 F_3 individuals, and selfing continued through to F_{10} . The focal individuals (i.e.
269 descendants that were the imputation targets) were F_2 , F_4 , F_6 , or F_{10} descendants.

270 In the base generation, individuals had their chromosomes sampled from the
271 100 base haplotypes. In subsequent generations the chromosomes of each individual
272 was sampled from parental chromosomes with recombination. The recombination rate
273 varied depending on the scenario resulting in chromosome sizes of 25, 50, 100, 150,
274 200, 300, and 400 centiMorgans (cM). Recombinations occurred with a 1%
275 probability per cM and were uniformly distributed along the chromosome.

276 *Simulated SNP marker arrays*

277 A single HD array of 25,000 SNP markers for the single chromosome was
278 simulated. To test the effect of the number of markers on the LD array, eight LD
279 arrays of 3, 5, 10, 20, 50, 100, 200, and 400 markers for the single chromosome were
280 simulated. Arrays were constructed by aiming to select a set of markers that
281 segregated in the parents and that were evenly distributed across the chromosome. All
282 LD arrays were nested within each other and within the HD array.

283 *Scenarios*

284 The imputation accuracy of AlphaPlantImpute and PlantImpute were
285 compared in four different scenarios (scenario 1, 2, 3, and 4). Scenarios 1, 2, and 3
286 were the same as scenarios 2, 4, and 5 in Hickey et al. 2015. A description of all four
287 scenarios is provided below. In all scenarios, focal individuals genotyped at LD were
288 imputed to the single HD array of 25,000 SNP markers. Ten replications of each
289 scenario were performed and the average of each replication is reported in the results.

290 Scenario 1: The effect of the number of selfing events separating parents and
291 focal individuals. Parents were almost fully inbred (F_{20}) and chromosomes were 100
292 cM in length. The accuracy of imputation was assessed for F_2 , F_4 , F_6 , and F_{10} focal
293 individuals.

294 Scenario 2: The effect of the level of inbreeding in parents. Parents were F_1 ,
295 F_2 , F_4 , F_{10} , or F_{20} and chromosomes were 100 cM in length. The accuracy of
296 imputation was assessed for F_2 focal individuals.

297 Scenario 3: The effect of chromosome size. Parents were fully inbred (F_{20}) and
298 the accuracy of imputation was assessed for F_2 focal individuals. Chromosomes were
299 25, 50, 100, 150, 200, 300, or 400 cM in size.

300 Scenario 4: The effect of number of focal individuals in the bi-parental
301 population. Parents were fully inbred (F_{20}) and the accuracy of imputation was
302 assessed for F_2 focal individuals. Subsets of focal individuals were randomly selected
303 from the 100 focal individuals to generate bi-parental population sizes of 1, 5, 10, 25,
304 and 50 focal individuals.

305 *Analysis*

306 Imputation was performed within each bi-parental population. Parents were
307 assumed genotyped at HD and focal individuals were assumed genotyped at LD. The
308 imputation accuracy was calculated for each focal individual as the correlation
309 between the true and imputed genotypes. The precision in imputation accuracy was
310 calculated as the log of the inverse of the variance in imputation accuracy within each
311 bi-parental population.

312

313 **Results**

314 For each scenario, we first present the imputation accuracy of
315 AlphaPlantImpute and then compare it to PlantImpute (Nettelblad et al., 2009; Hickey
316 et al., 2015).

317 *Effect of the number of markers on the low-density array*

318 Increasing the number of LD markers increases the imputation accuracy of
319 AlphaPlantImpute. Figure 2 plots the number of LD markers against the accuracy of
320 imputation for F_2 focal individuals of an $F_{20} \times F_{20}$ bi-parental cross. Figure 2 shows
321 that increasing the number of LD markers from 3 to 20 SNP increased the average
322 imputation accuracy from 0.85 to 0.96. Increasing the number of markers beyond 20
323 achieved only a slight increase in the accuracy of imputation from 0.96 with 20
324 markers to >0.99 with 400 markers.

325 *Scenario 1: Effect of the number of selfing events separating parents and focal* 326 *individuals*

327 Increasing the number of selfing events separating parents and focal
328 individuals slightly decreases the imputation accuracy of AlphaPlantImpute. Figure
329 3a plots the accuracy of imputation in F_2 , F_4 , F_6 and F_{10} focal individuals of a bi-
330 parental population where the parents were F_{20} . Figure 3a shows that with 3 LD
331 markers, the average imputation accuracy decreased from 0.85 for F_2 focal individuals
332 to 0.77 for F_{10} focal individuals. Increasing the number of LD markers beyond 10
333 markers mitigates the decrease in the average imputation accuracy between F_2 focal
334 individuals and F_{10} focal individuals. Figure 3a shows that with 20 LD markers, the

335 average imputation accuracy decreased from 0.96 for F_2 focal individuals to 0.95 for
336 F_{10} focal individuals.

337 Regardless of the number of selfing events separating parents and focal
338 individuals, the accuracy of imputation for AlphaPlantImpute was higher than for
339 PlantImpute when the number of LD markers was low. Figure 3b plots the average
340 imputation accuracy of AlphaPlantImpute on the y-axis and for PlantImpute on the x-
341 axis. The colours represent the different number of LD markers and the shapes
342 represent the number of selfing events separating the parents and the focal
343 individuals. The red diagonal line indicates when the imputation accuracy of the two
344 methods is equal. Points above the line indicate when the accuracy of imputation was
345 higher for AlphaPlantImpute than for PlantImpute and visa versa. Figure 3b shows
346 that with 3 LD markers, the average accuracy of imputation was 0.85 for
347 AlphaPlantImpute and 0.76 for PlantImpute for F_2 focal individuals and was 0.77 for
348 AlphaPlantImpute and 0.70 for PlantImpute for F_{10} focal individuals.

349 For all numbers of selfing events separating parents and focal individuals,
350 increasing the number of LD markers reduced and in some cases reversed the
351 advantage of AlphaPlantImpute over PlantImpute. This was most obvious for F_{10}
352 focal individuals for medium number of LD markers where the imputation accuracy
353 with PlantImpute was slightly higher than with AlphaPlantImpute. Figure 3b shows
354 that with 10 LD markers, the average imputation accuracy was 0.93 for
355 AlphaPlantImpute and 0.94 for PlantImpute for F_2 focal individuals and was 0.90 for
356 AlphaPlantImpute and 0.92 for PlantImpute for F_{10} focal individuals. Increasing the
357 number of LD markers beyond 100 markers meant that the average accuracy of
358 imputation for AlphaPlantImpute equalled that for PlantImpute. Figure 3b shows that

359 with 100 LD markers, the average imputation accuracy was 0.99 for both
360 AlphaPlantImpute and PlantImpute for F_2 focal individuals and for F_{10} focal
361 individuals.

362 For all numbers of selfing events separating parents and focal individuals, the
363 precision of imputation accuracy (i.e., consistency across focal individuals) for
364 AlphaPlantImpute was higher than for PlantImpute when the number of LD markers
365 was low. Figure 3c is similar to Figure 3b and plots the log of the precision of
366 imputation accuracy for AlphaPlantImpute on the y-axis and PlantImpute on the x-
367 axis. Points above the line indicate better precision (i.e. less variance) for
368 AlphaPlantImpute than for PlantImpute, and vice versa. Figure 3c shows that with 3
369 LD markers, the precision of imputation was 1.62 for AlphaPlantImpute and 1.08 for
370 PlantImpute for F_2 focal individuals and was 1.32 for AlphaPlantImpute and 1.11 for
371 PlantImpute for F_{10} focal individuals.

372 Figure 3c also shows that for medium number of LD markers, the precision of
373 imputation accuracy for AlphaPlantImpute was higher than for PlantImpute for F_2
374 focal individuals but was lower when the number of selfing events was higher. With
375 20 LD markers, the precision of imputation accuracy was 2.48 for AlphaPlantImpute
376 and 2.00 for PlantImpute for F_2 focal individuals and was 2.57 for AlphaPlantImpute
377 and 2.80 for PlantImpute for F_{10} focal individuals. With the highest number of LD
378 markers (400), the precision of imputation accuracy was 3.84 for AlphaPlantImpute
379 and 4.00 for PlantImpute for F_2 focal individuals and was 5.40 for both
380 AlphaPlantImpute and PlantImpute for F_{10} focal individuals.

381 *Scenario 2: Effect of the level of inbreeding in parents*

382 Increasing the level of inbreeding in the parents increases the imputation
383 accuracy for AlphaPlantImpute. Figure 4a plots the accuracy of imputation in F_2 focal
384 individuals of a bi-parental population where the parents were F_1 , F_2 , F_4 , F_{10} or F_{20} .
385 Figure 4a shows that with 20 LD markers, the average imputation accuracy increased
386 from 0.81 for F_1 parents to 0.96 for F_{20} parents. Figure 4a also shows that increasing
387 the level of inbreeding in the parents beyond F_4 did not increase the average accuracy
388 of imputation for F_2 focal individuals. The average imputation accuracy with 20 LD
389 markers was approximately 0.96 for F_2 focal individuals when parents were F_4 , F_{10} ,
390 and F_{20} .

391 For all levels of inbreeding in the parents and all numbers of LD markers, the
392 average imputation accuracy with AlphaPlantImpute was almost always higher than
393 with PlantImpute. Figure 4b is similar to Figure 3b and plots the average imputation
394 accuracy for AlphaPlantImpute on the y-axis and for PlantImpute on the x-axis. The
395 shapes represent the level of inbreeding in the parents. Figure 4b shows that with 20
396 SNP LD markers, the average imputation accuracy was 0.81 for AlphaPlantImpute
397 and 0.74 for PlantImpute for F_2 focal individuals when parents were F_1 , 0.95 for
398 AlphaPlantImpute and 0.91 for PlantImpute when parents were F_4 , and 0.96 for
399 AlphaPlantImpute and 0.94 for PlantImpute when parents were F_{10} . In two cases, the
400 average imputation accuracy with PlantImpute was slightly higher than with
401 AlphaPlantImpute. This was when parents were F_4 and with 3 and 5 LD markers. The
402 average imputation accuracy was 0.84 for AlphaPlantImpute and 0.80 for PlantImpute
403 with 3 LD markers and was 0.87 for AlphaPlantImpute and 0.85 for PlantImpute with
404 5 LD markers.

405 For all levels of inbreeding in the parents and all numbers of LD markers, the
406 precision of imputation accuracy with AlphaPlantImpute was almost always higher
407 than with PlantImpute. Figure 4c is similar to 3c and plots the log of the precision of
408 imputation accuracy for AlphaPlantImpute on the y-axis and PlantImpute on the x-
409 axis. Figure 4c shows that with 20 LD markers, the precision of imputation accuracy
410 was 2.16 for AlphaPlantImpute and 1.92 for PlantImpute for F_2 focal individuals
411 when parents were F_1 , 2.54 for AlphaPlantImpute and 1.84 for PlantImpute when
412 parents were F_4 , and 2.52 for AlphaPlantImpute and 1.71 for PlantImpute when
413 parents were F_{10} . In a few cases, the precision of imputation accuracy for PlantImpute
414 was slightly higher than AlphaPlantImpute. This was mainly when parents were F_{20}
415 and with 50, 200, and 400 LD markers. The precision of imputation accuracy was
416 3.04 for AlphaPlantImpute and 3.40 for PlantImpute with 50 LD markers, was 3.71
417 for AlphaPlantImpute and 4.00 for PlantImpute with 200 LD markers, and was 3.84
418 for AlphaPlantImpute and 4.00 for PlantImpute with 400 LD markers.

419 *Scenario 3: Effect of chromosome size*

420 Increasing the chromosome size (in cM) decreased the imputation accuracy
421 for AlphaPlantImpute. This was most apparent when the number of LD markers was
422 10 or less. Figure 5a plots the imputation accuracy for seven chromosome sizes of 25,
423 50, 100, 150, 200, 300, and 400 cM for F_2 focal individuals of a bi-parental
424 population where the parents were F_{20} . Figure 5a shows that with 3 LD markers,
425 quadrupling the chromosome size from 25 cM to 100 cM decreased the average
426 imputation accuracy from 0.95 to 0.85, and quadrupling from 100 cM to 400 cM
427 decreased the average imputation accuracy from 0.85 to 0.55. The reduction in the
428 imputation accuracy was less or non-existent when the number of LD markers was

429 higher than 10. Figure 5a shows that the imputation accuracy was approximately 0.98
430 for all chromosome sizes when the number of LD markers was 50.

431 When the chromosome size was 300 cM or less, the average imputation
432 accuracy was higher for AlphaPlantImpute than for PlantImpute. Figure 5b is similar
433 to Figure 3b and plots the average imputation accuracy for AlphaPlantImpute on the
434 y-axis and for PlantImpute on the x-axis. The shapes represent the chromosome sizes.
435 Figure 5b shows that with 3 LD markers, the average imputation accuracy was 0.95
436 for AlphaPlantImpute and 0.69 for PlantImpute when the chromosome size was 25
437 cM and was 0.61 for AlphaPlantImpute and 0.57 for PlantImpute when the
438 chromosome size was 300 cM. The exception to this was when the chromosome size
439 was 150 cM, where the average imputation accuracy was 0.70 for AlphaPlantImpute
440 and 0.83 for PlantImpute. When the chromosome size was 400 cM the average
441 imputation accuracy was 0.55 for AlphaPlantImpute and 0.51 for PlantImpute when 3
442 LD markers were used but was 0.61 for AlphaPlantImpute and 0.68 for PlantImpute
443 when 5 LD markers were used.

444 For all chromosome sizes and numbers of LD markers, the precision of
445 imputation accuracy for AlphaPlantImpute was generally higher than for PlantImpute.
446 Figure 5c is similar to Figure 3c and plots the precision of imputation accuracy for
447 AlphaPlantImpute on the y-axis and for PlantImpute on the x-axis. Figure 5c shows
448 that with 3 LD markers, the precision of imputation accuracy was 0.71 for
449 AlphaPlantImpute and 1.78 for PlantImpute when the chromosome size was 25 cM,
450 was 1.08 for AlphaPlantImpute and 1.62 for PlantImpute when the chromosome size
451 was 100 cM and was 1.59 for AlphaPlantImpute and 1.20 for PlantImpute when the
452 chromosome size was 400 cM. The exception to this was when the chromosome size

453 was 150 cM, where the precision of imputation accuracy was 1.17 for
454 AlphaPlantImpute and 1.46 for PlantImpute.

455 *Scenario 4: Effect of the number of focal individuals in the bi-parental population*

456 Increasing the number of focal individuals in the bi-parental population
457 slightly increased the imputation accuracy for AlphaPlantImpute. This was most
458 apparent when the number of LD markers was low. Figure 6 plots the accuracy of
459 imputation for F_2 focal individuals of an $F_{20} \times F_{20}$ bi-parental cross with 1, 5, 10, 25,
460 50 or 100 focal individuals. Figure 6 shows that increasing the number of focal
461 individuals from 5 to 100 increased the average imputation accuracy from 0.83 to
462 0.85 when 3 LD markers were used. Figure 6 also shows that when the 10 or more LD
463 markers were used, increasing the number of focal individuals had no effect on the
464 imputation accuracy. When the number of LD markers was 400, the average
465 imputation accuracy was 0.96 with 5 or 100 focal individuals in the bi-parental
466 population.

467 Figure 6 also shows that when we only imputed one focal individual, the
468 imputation accuracy fluctuated according to the focal individual that was sampled. As
469 a result, increasing the number of LD markers did not always increase the imputation
470 accuracy. For example, the average imputation accuracy was 0.95, 0.91, or 0.94 when
471 3, 5, or 10 LD markers were used. When 400 LD markers were used, the average
472 accuracy of imputation was 0.997.

473 *Computational requirements of AlphaPlantImpute*

474 Table 1 summarises the computational requirements of AlphaPlantImpute for
475 twelve datasets across the three scenarios. Datasets were chosen to reflect the

476 extremes in the number of selfing events separating parents and focal individuals (F_2
477 vs. F_{10}), the level of inbreeding in the parents (F_1 vs. F_{20}) and the number of LD
478 markers (3, 50, or 400). Table 1 shows that the average run time for
479 AlphaPlantImpute was 22.13 seconds with a maximum of 49.33 seconds. The average
480 memory requirement for AlphaPlantImpute was 0.08 GB with a maximum of 0.082
481 GB.

482

483 **Discussion**

484 Our results highlight three points for discussion: (i) the performance of
485 AlphaPlantImpute; (ii) the performance of AlphaPlantImpute compared to
486 PlantImpute; and (iii) future development of AlphaPlantImpute.

487 *Performance of AlphaPlantImpute*

488 This paper presents a new heuristic method, called AlphaPlantImpute, for
489 phasing and imputation of SNP array data in diploid plant species. AlphaPlantImpute
490 explicitly leverages features of plant breeding programs to impute LD focal
491 individuals to HD. The explicit utilisation of pedigree information and heuristics
492 developed specifically to track the inheritance of parental haplotypes using the LD
493 genotypes of focal individuals are likely to be the reasons for AlphaPlantImpute's
494 robust and consistent performance across all tested scenarios. AlphaPlantImpute
495 achieves high imputation accuracy of between 0.8 and 1.0 for the majority of
496 scenarios. For scenarios where the imputation accuracy was below 0.8, increasing the
497 number of LD markers increased the imputation accuracy.

498 Increasing number of selfing events separating parents and focal individuals
499 from F_2 to F_{10} only slightly decreases the imputation accuracy. Decreasing the level of
500 inbreeding in the parents or increasing the chromosome size decreases the imputation
501 accuracy when the number of LD markers is 10 or less. However, in both cases, the
502 decrease in the imputation accuracy could be mitigated by increasing the number of
503 LD markers to 20 SNP or more.

504 Decreasing the number of focal individuals in the bi-parental population
505 slightly decreases the imputation accuracy. This was most evident when the number

506 of LD markers was 10 SNP or less. The likely cause of this is that inferring the most
507 likely linkage between alleles for two markers is difficult with fewer focal
508 individuals, since fewer individuals will be homozygous at the markers. In this case,
509 the algorithm defaults to the linkage pattern of alleles in the parents. This may be sub-
510 optimal for imputing markers in regions with elevated recombination rates, i.e.,
511 hotspots. When there was a single focal individual in the focal family, the accuracy of
512 imputation for that individual varied. The likely cause of this is whether an individual
513 had a recombination or whether it had inherited the parental haplotypes without
514 recombination. One solution to this situation could be to utilise the most likely
515 linkage from related families with more genotyped focal individuals (see section:
516 Future work and developments).

517 Overall, the results suggest that for a given population, high imputation
518 accuracy can be achieved even when the number of LD markers is low, and small
519 increases in the number of markers can achieve high accuracies depending on the
520 biology of the species (i.e. recombination rate, obligate outcrossing) and the pedigree
521 design (outbred, inbred, level of selfing).

522 *Performance of AlphaPlantImpute compared to PlantImpute*

523 The imputation accuracy for AlphaPlantImpute was compared to that for
524 PlantImpute (Nettelblad et al., 2009; Hickey et al., 2015). In the majority of cases, the
525 imputation accuracy was higher for AlphaPlantImpute than for PlantImpute. One
526 exception to this was when the chromosome size was 400 cM and when the number
527 of LD markers was 20 or less (e.g. 0.88 vs. 0.90 when the number of LD markers was
528 20). One reason for this could be that unless there is enough information in the
529 genotypes of focal individual on the LD array, the heuristic algorithm in

530 AlphaPlantImpute is inherently more conservative in determining recombination
531 regions compared to the probabilistic algorithm in PlantImpute. As such,
532 AlphaPlantImpute is more likely to leave positions as missing and fill them in as the
533 parent average in the final step.

534 The precision of imputation accuracy (calculated as the log of the inverse of
535 the variance in imputation accuracy within each bi-parental population) was also
536 higher in the majority of cases for AlphaPlantImpute than for PlantImpute. This was
537 most apparent with small number of LD markers. The higher precision of imputation
538 accuracy for AlphaPlantImpute is likely a consequence of directly calling allele phase
539 and parent-of-origin and imputed genotypes in turn. The probabilistic algorithm of
540 PlantImpute is marginalizing over the all possible phase and genotype, which is
541 probabilistically correct and handles the uncertainty properly, but it seems this is
542 lowering the imputation accuracy. One exception to this was when the chromosome
543 size was 150 cM, where the precision of imputation accuracy was higher for
544 PlantImpute than for AlphaPlantImpute for all LD arrays.

545 The biggest advantage of AlphaPlantImpute compared to PlantImpute relates
546 to computational requirements. Hickey et. al. 2015 report that to perform imputation
547 within a single bi-parental population of 100 F₂ focal individuals, PlantImpute
548 required a minimum of 3 hours and in excess of 100 GB of memory. In comparison,
549 AlphaPlantImpute required on average ~22 seconds and ~0.08 GB of memory for all
550 tested scenarios.

551 The high and consistent accuracies achieved with very low computational
552 requirements makes AlphaPlantImpute an attractive, reliable and practical tool for

553 routine use in plant breeding programs that are already using or will include SNP
554 array data to inform selection decisions.

555 *Future work and developments*

556 At present, the heuristic method in AlphaPlantImpute works within the most
557 common plant breeding program design of bi-parental populations and it works best
558 when parents are fully inbred or close to being fully inbred. AlphaPlantImpute could
559 be extended in multiple ways. For example, instead of treating each bi-parental
560 population as an independent unit it could simultaneously work across bi-parental
561 populations that share parents. This could increase the imputation accuracy in three
562 ways: (i) information between bi-parental populations could be shared for imputation
563 of focal individuals that are effectively half-sibs (one common parent); (ii)
564 information between bi-parental populations could be used to resolve phase where
565 one or both parents are heterozygous at one or more consecutive markers; and (iii) if a
566 common parent has no or LD genotypes available, information from its descendants
567 across half-sib bi-parental populations could be leveraged to phase and impute it to
568 high-density.

569 AlphaPlantImpute could also be extended to include ancestral pedigree
570 information (such as grandparents and great-grandparents). This could be useful for
571 improving phasing and imputation of parents with missing information or that are
572 highly outbred. More simply, AlphaPlantImpute could also be extended so that it can
573 directly read in and exploit phased information for the fully or partially outbred
574 parents. Such phased information could be generated for parents by running
575 AlphaPlantImpute on the bi-parental family from which the fully or partially outbred
576 parent derived.

577 AlphaPlantImpute could be extended so that it reads in previously inferred
578 most likely linked alleles at two markers. It is likely that linkage patterns are shared
579 across families, especially if the families are related. Using this information across
580 families would be especially suited to imputation situations in bi-parental populations
581 that have only a few genotyped focal individuals (e.g., one genotyped individual per
582 family).

583 Finally, although SNP arrays for the many domesticated plant species exist,
584 low-coverage sequencing methods such as genotyping-by-sequencing are also used.
585 The heuristics of AlphaPlantImpute might be extended to enable imputation with such
586 data.

587 *Software availability*

588 We implemented our method in a software package called AlphaPlantImpute,
589 which is available for download at
590 <http://www.AlphaGenes.roslin.ed.ac.uk/AlphaPlantImpute/> along with a user manual.

591

592 **Acknowledgments**

593 The authors acknowledge the financial support from the BBSRC ISP grant

594 number ‘BB/P013759/1’ and from the BBSRC KWS grant number ‘BB/R002061/1’.

595 This work has made use of the resources provided by the Edinburgh Compute and

596 Data Facility (ECDF) (<http://www.ecdf.ed.ac.uk>).

597

598

599 **References**

- 600 Antolín, R., C. Nettelblad, G. Gorjanc, D. Money, and J.M. Hickey. 2017. A hybrid
601 method for the imputation of genomic data in livestock populations. *Genet.*
602 *Sel. Evol.* 49(1): 30. doi: 10.1186/s12711-017-0300-y.
- 603 Chen, G.K., P. Marjoram, and J.D. Wall. 2009. Fast and flexible simulation of DNA
604 sequence data. *Genome Res.* 19(1): 136–142. doi: 10.1101/gr.083634.108.
- 605 Cleveland, M.A., and J.M. Hickey. 2013. Practical implementation of cost-effective
606 genomic selection in commercial pig breeding using imputation. *J. Anim. Sci.*
607 91(8): 3583–3592. doi: 10.2527/jas.2013-6270.
- 608 Druet, T., and M. Georges. 2010. A Hidden Markov Model Combining Linkage and
609 Linkage Disequilibrium Information for Haplotype Reconstruction and
610 Quantitative Trait Locus Fine Mapping. *Genetics* 184(3): 789–798. doi:
611 10.1534/genetics.109.108431.
- 612 Faux, A.-M., G. Gorjanc, R.C. Gaynor, M. Battagin, S.M. Edwards, D.L. Wilson, S.J.
613 Hearne, S. Gonen, and J.M. Hickey. 2016. AlphaSim: Software for Breeding
614 Program Simulation. *Plant Genome* 9(3). doi:
615 10.3835/plantgenome2016.02.0013.
- 616 Gorjanc, G., M. Battagin, J.-F. Dumasy, R. Antolin, R.C. Gaynor, and J.M. Hickey.
617 2017a. Prospects for Cost-Effective Genomic Selection via Accurate Within-
618 Family Imputation. *Crop Sci.* 57(1): 216. doi: 10.2135/cropsci2016.06.0526.
- 619 Gorjanc, G., J.-F. Dumasy, S. Gonen, R.C. Gaynor, R. Antolin, and J.M. Hickey.
620 2017b. Potential of Low-Coverage Genotyping-by-Sequencing and Imputation
621 for Cost-Effective Genomic Selection in Biparental Segregating Populations.
622 *Crop Sci.* 57(3): 1404–1420. doi: 10.2135/cropsci2016.08.0675.
- 623 Hickey, J.M., G. Gorjanc, R.K. Varshney, and C. Nettelblad. 2015. Imputation of
624 Single Nucleotide Polymorphism Genotypes in Biparental, Backcross, and
625 Topcross Populations with a Hidden Markov Model. *Crop Sci.* 55(5): 1934–
626 1946. doi: 10.2135/cropsci2014.09.0648.
- 627 Hickey, J.M., B.P. Kinghorn, B. Tier, J.F. Wilson, N. Dunstan, and J.H. van der Werf.
628 2011. A combined long-range phasing and long haplotype imputation method
629 to impute phase for SNP genotypes. *Genet. Sel. Evol.* 43(1): 12. doi:
630 10.1186/1297-9686-43-12.
- 631 Hickey, J.M., and A. Kranis. 2013. Extending long-range phasing and haplotype
632 library imputation methods to impute genotypes on sex chromosomes. *Genet.*
633 *Sel. Evol.* 45(1): 10. doi: 10.1186/1297-9686-45-10.
- 634 Howie, B.N., P. Donnelly, and J. Marchini. 2009. A flexible and accurate genotype
635 imputation method for the next generation of genome-wide association
636 studies. *PLoS Genet.* 5(6): e1000529.

- 637 Jacobson, A., L. Lian, S. Zhong, and R. Bernardo. 2014. General Combining Ability
638 Model for Genomewide Selection in a Biparental Cross. *Crop Sci.* 54(3): 895.
639 doi: 10.2135/cropsci2013.11.0774.
- 640 Jacobson, A., L. Lian, S. Zhong, and R. Bernardo. 2015. Marker imputation before
641 genomewide selection in biparental maize populations. *Plant Genome* 8(2): 9.
642 doi: doi:10.3835/plantgenome2014.10.0078.
- 643 Kong, A., G. Masson, M.L. Frigge, A. Gylfason, P. Zusmanovich, G. Thorleifsson,
644 P.I. Olason, A. Ingason, S. Steinberg, T. Rafnar, P. Sulem, M. Mouy, F.
645 Jonsson, U. Thorsteinsdottir, D.F. Gudbjartsson, H. Stefansson, and K.
646 Stefansson. 2008. Detection of sharing by descent, long-range phasing and
647 haplotype imputation. *Nat. Genet.* 40(9): 1068–1075. doi: 10.1038/ng.216.
- 648 Li, Y., C.J. Willer, J. Ding, P. Scheet, and G.R. Abecasis. 2010. MaCH: using
649 sequence and genotype data to estimate haplotypes and unobserved genotypes.
650 *Genet. Epidemiol.* 34(8): 816–834. doi: 10.1002/gepi.20533.
- 651 Loh, P.-R., P. Danecek, P.F. Palamara, C. Fuchsberger, Y. A Reshef, H. K Finucane,
652 S. Schoenherr, L. Forer, S. McCarthy, G.R. Abecasis, R. Durbin, and A. L
653 Price. 2016. Reference-based phasing using the Haplotype Reference
654 Consortium panel. *Nat. Genet.* 48(11): 1443–1448. doi: 10.1038/ng.3679.
- 655 Nettelblad, C., S. Holmgren, L. Crooks, and Ö. Carlborg. 2009. cnF2freq: Efficient
656 Determination of Genotype and Haplotype Probabilities in Outbred
657 Populations Using Markov Models. p. 307–319. *In* Rajasekaran, S. (ed.),
658 *Bioinformatics and Computational Biology. Lecture Notes in Computer
659 Science.* Springer Berlin Heidelberg.
- 660 O’Connell, J., K. Sharp, N. Shrine, L. Wain, I. Hall, M. Tobin, J.-F. Zagury, O.
661 Delaneau, and J. Marchini. 2016. Haplotype estimation for biobank-scale data
662 sets. *Nat. Genet.* advance online publication. doi: 10.1038/ng.3583.
- 663 Sargolzaei, M., J.P. Chesnais, and F.S. Schenkel. 2011. FImpute - An efficient
664 imputation algorithm for dairy cattle populations. *J. Dairy Sci.* 94 (E-Suppl.
665 1): 421.
- 666 Swarts, K., H. Li, J.A. Romero Navarro, D. An, M.C. Romay, S. Hearne, C. Acharya,
667 J.C. Glaubitz, S. Mitchell, R.J. Elshire, E.S. Buckler, and P.J. Bradbury. 2014.
668 Novel Methods to Optimize Genotypic Imputation for Low-Coverage, Next-
669 Generation Sequence Data in Crop Plants. *Plant Genome* 7(3): 0. doi:
670 10.3835/plantgenome2014.05.0023.
- 671 VanRaden, P.M., C. Sun, and J.R. O’Connell. 2015. Fast imputation using medium or
672 low-coverage sequence data. *BMC Genet.* 16(1): 82. doi: 10.1186/s12863-
673 015-0243-7.

674

675

676 **Figure captions**

677 **Figure 1. Schematic of heuristic algorithm of AlphaPlantImpute.**

678 **Figure 2. Effect of the number of SNP on the low-density array.**

679 **Figure 3. Effect of level of inbreeding in focal individuals.**

680 **Figure 4. Effect of the level of inbreeding in parents.**

681 **Figure 5. Effect of chromosome size.**

682

683 **Table captions**

684 **Table 1. Computational requirements of AlphaPlantImpute.**

685

686 **Supplementary Files**

687 **Supplementary File 1. Detailed schematic of heuristic algorithm of**
688 **AlphaPlantImpute.**

689

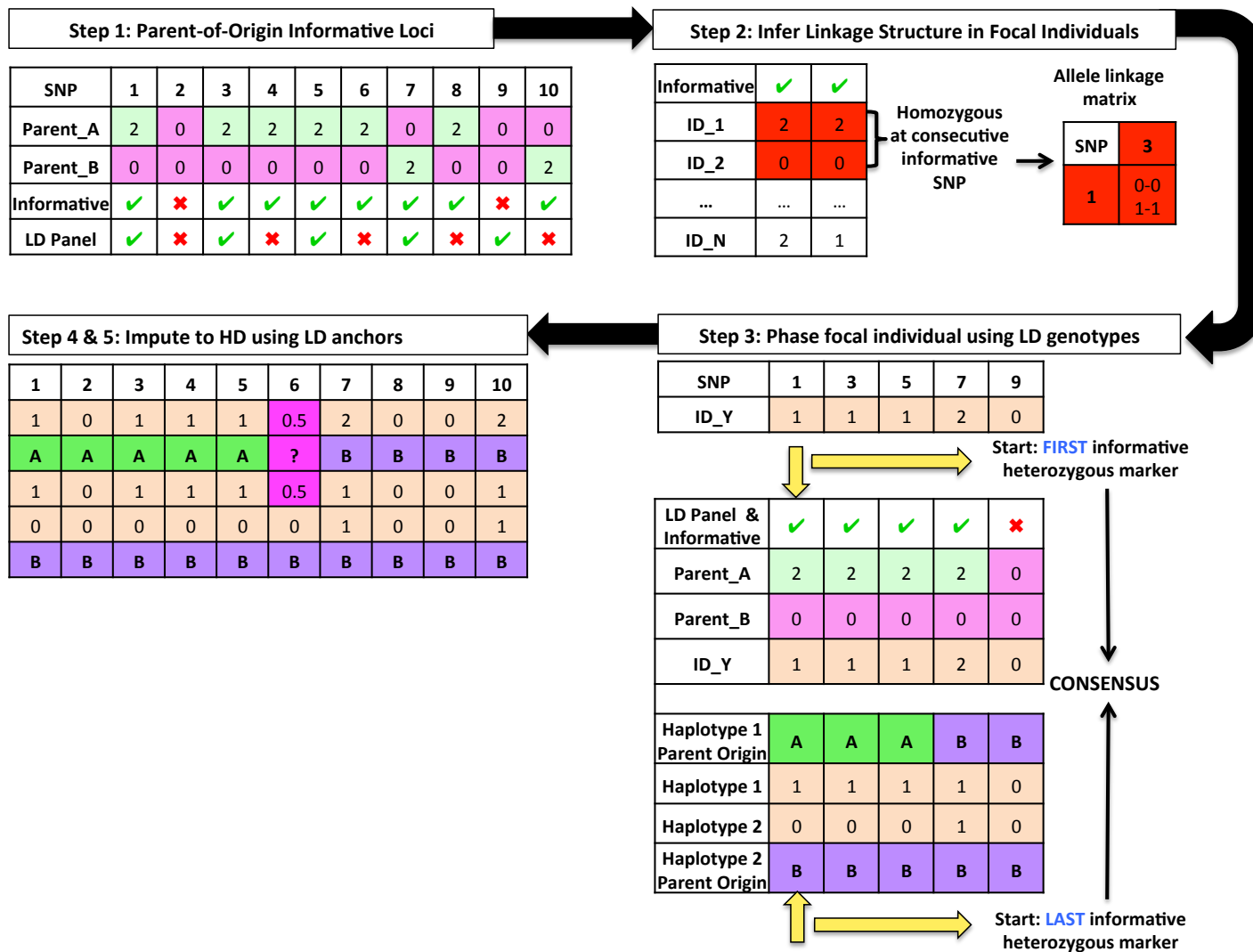


Figure 1 – Schematic of heuristic algorithm of AlphaPlantImpute

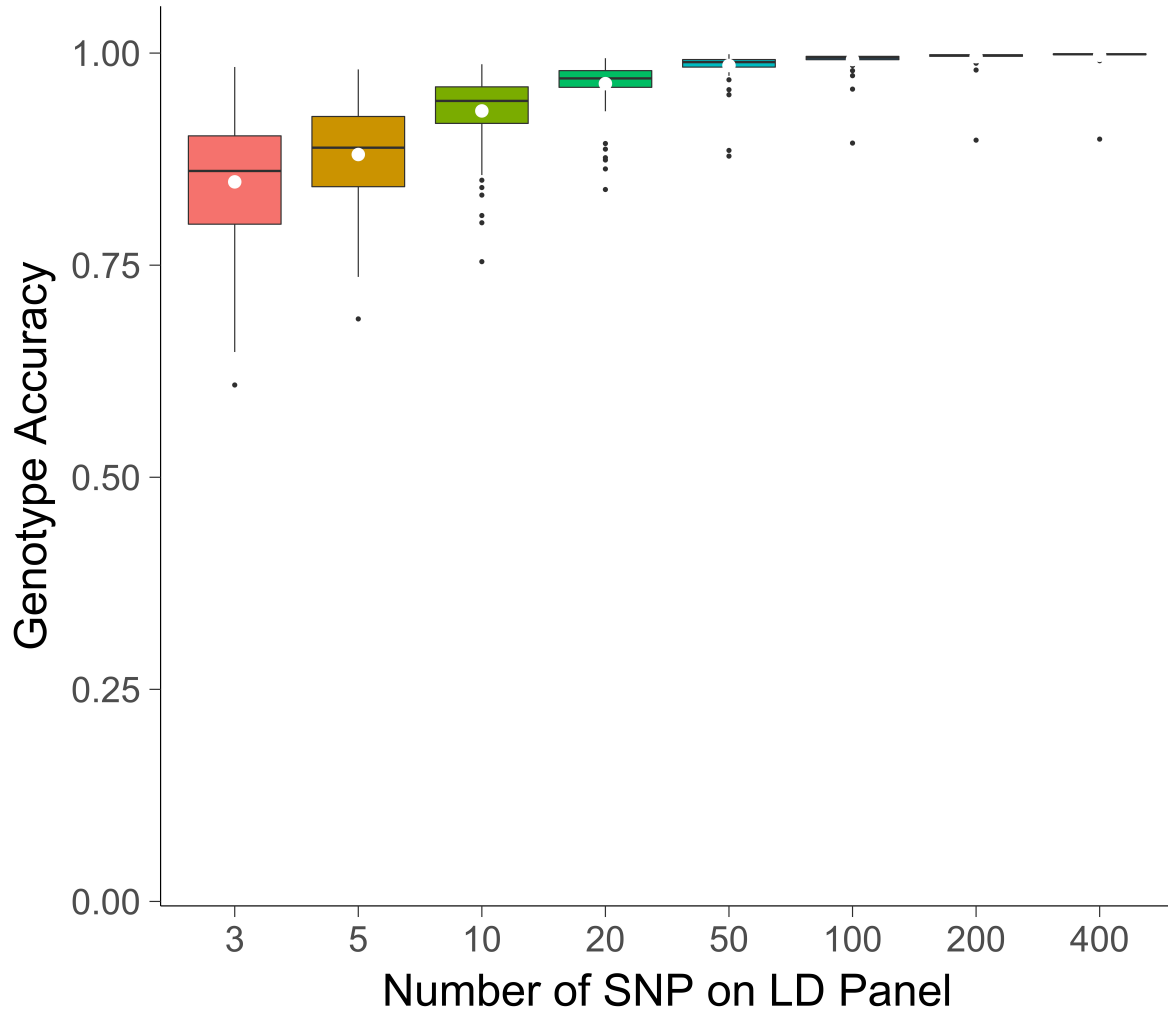
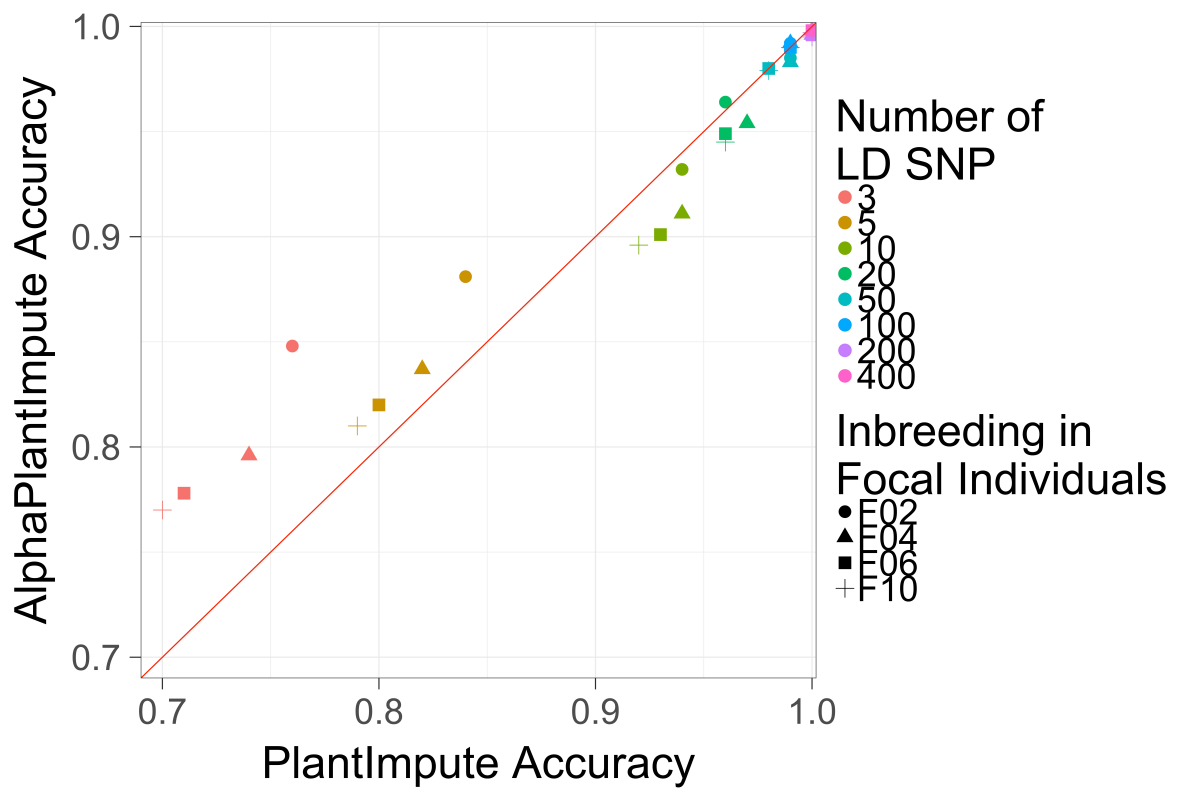
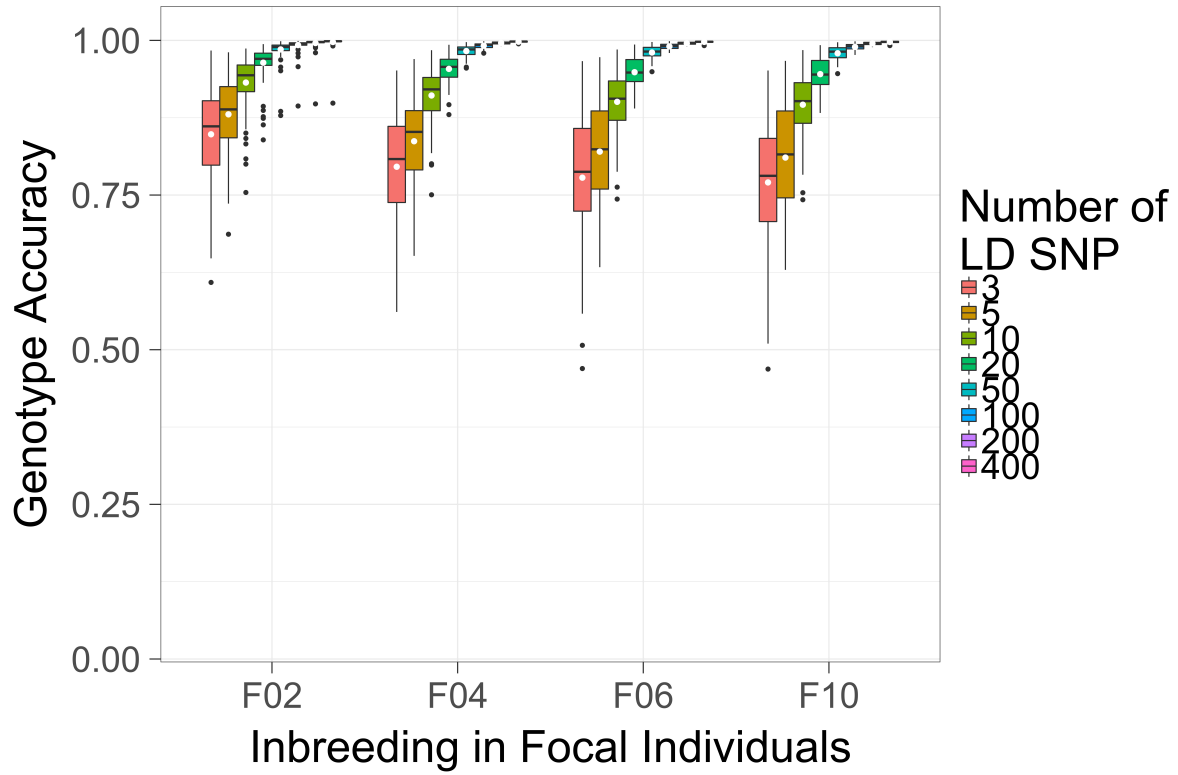


Figure 2 – Effect of the number of SNP on the low-density panel.

The number of SNP on the LD panel against the genotype imputation accuracy using AlphaPlantImpute for F_2 focal individuals of a bi-parental cross where the parents are F_{20} inbred individuals.



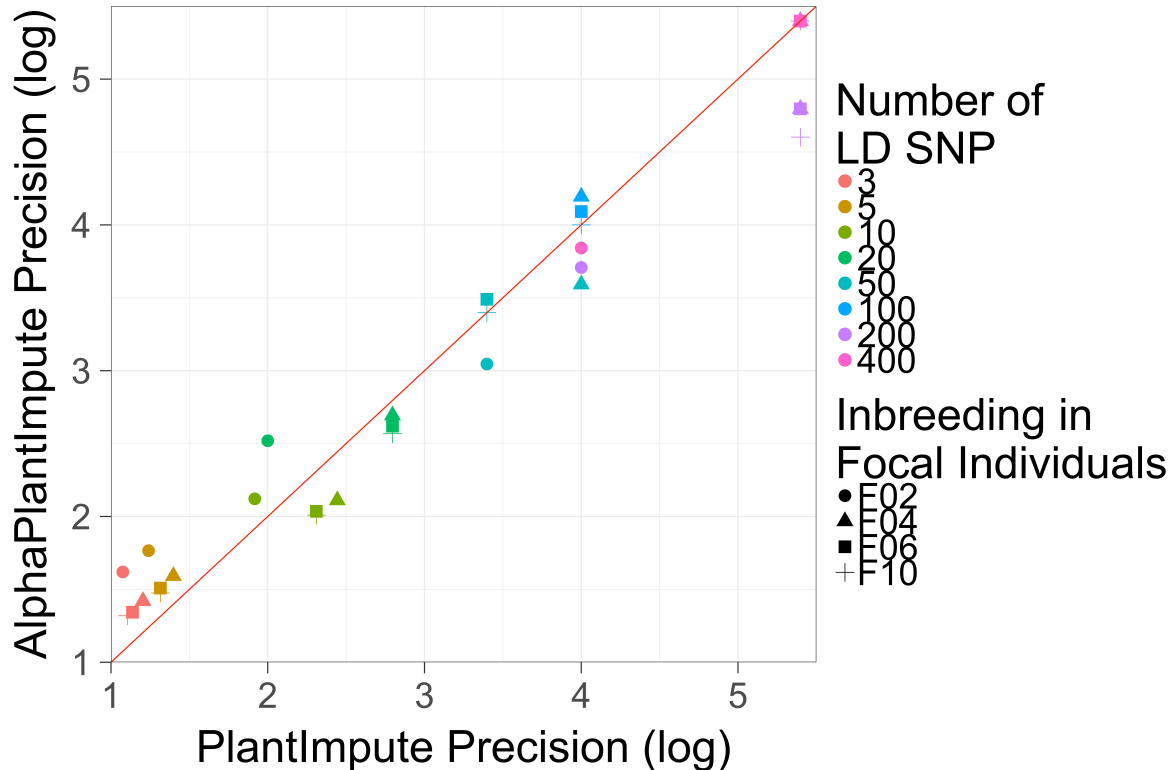
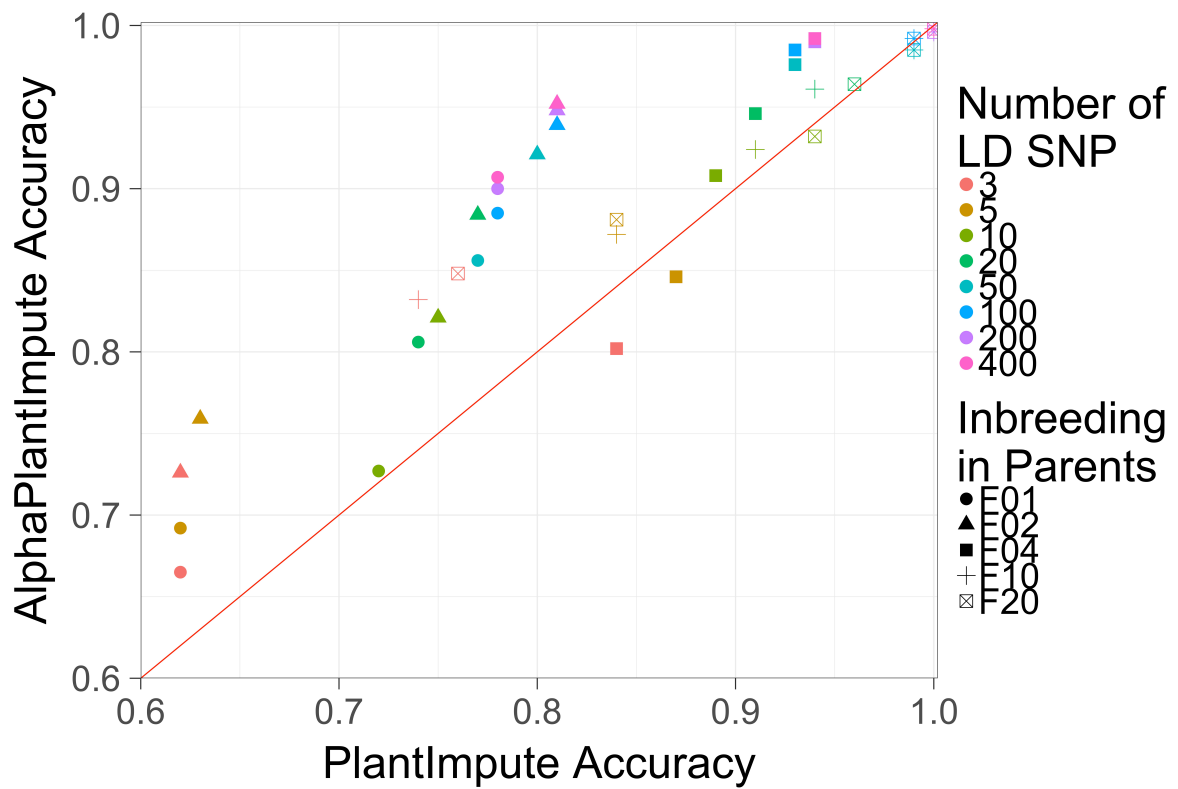
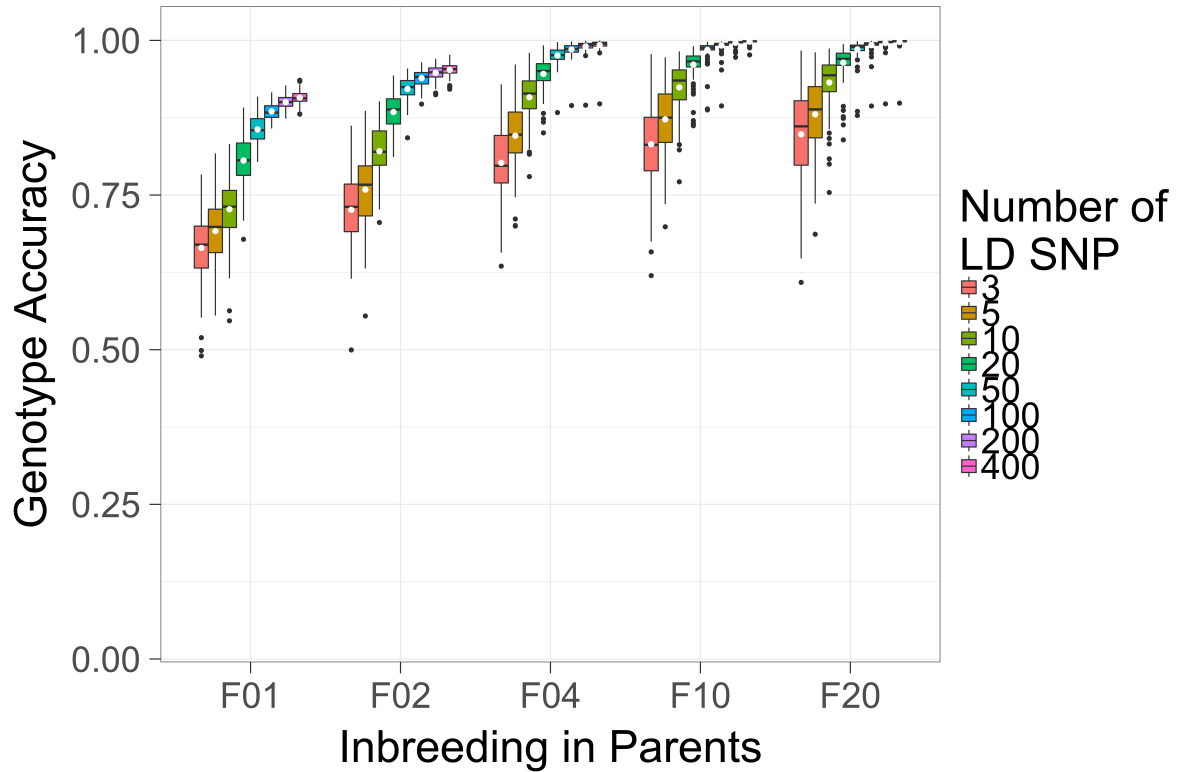


Figure 3 – Effect of the level of inbreeding in focal individuals.

(a) The genotype imputation accuracy using AlphaPlantImpute in F₂, F₄, F₆ and F₁₀ focal individuals from a bi-parental cross where the parents are F₂₀ inbred individuals.

(b) Comparison of the average genotype imputation accuracy using AlphaPlantImpute (y-axis) vs. PlantImpute (x-axis). The colours represent the different LD panels. The shapes represent the level of inbreeding in the focal individuals. The red diagonal line indicates when the accuracy of PlantImpute equals AlphaPlantImpute. Points above the line are when imputation accuracy is higher with AlphaPlantImpute and points below the line are when imputation accuracy is higher with PlantImpute.

(c) Comparison of the precision in imputation accuracy using AlphaPlantImpute (y-axis) vs. using PlantImpute (x-axis). The colours represent the different LD panels. The shapes represent the level of inbreeding in the focal individuals. The red diagonal line indicates when the precision of PlantImpute equals AlphaPlantImpute. Points above the line indicate when the precision in accuracies is higher in AlphaPlantImpute and points below the line are when the precision in accuracies is higher in PlantImpute.



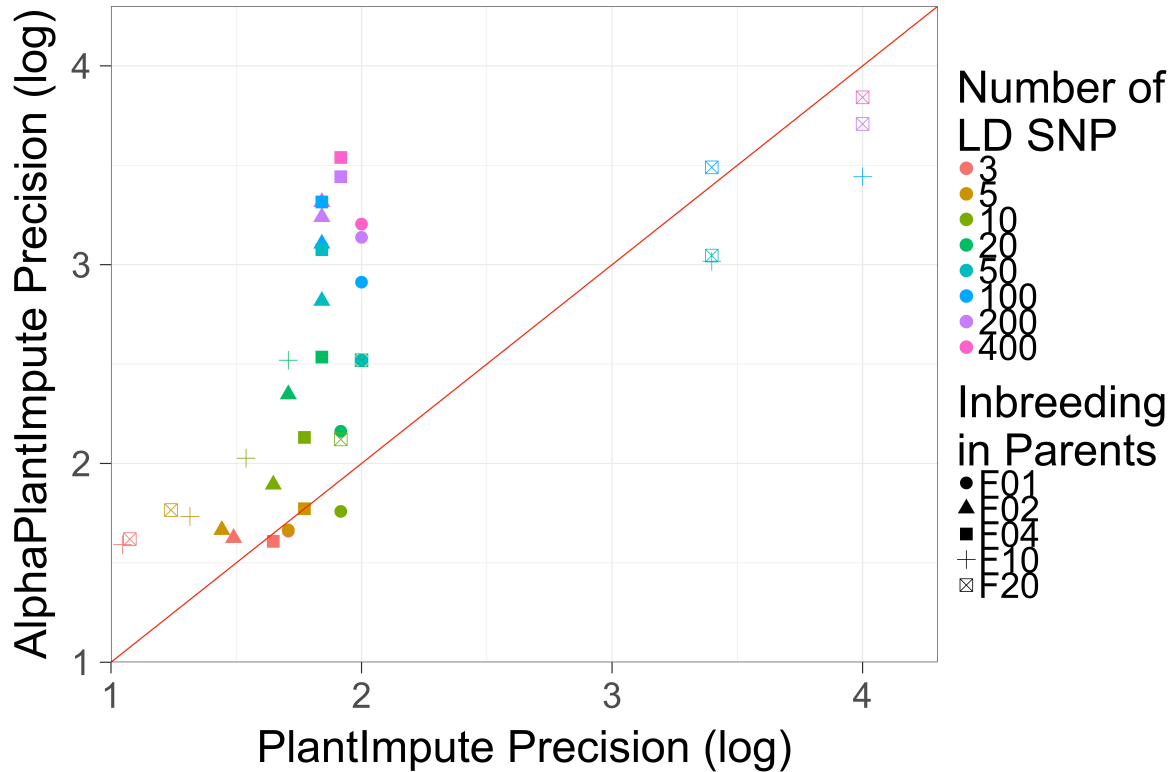
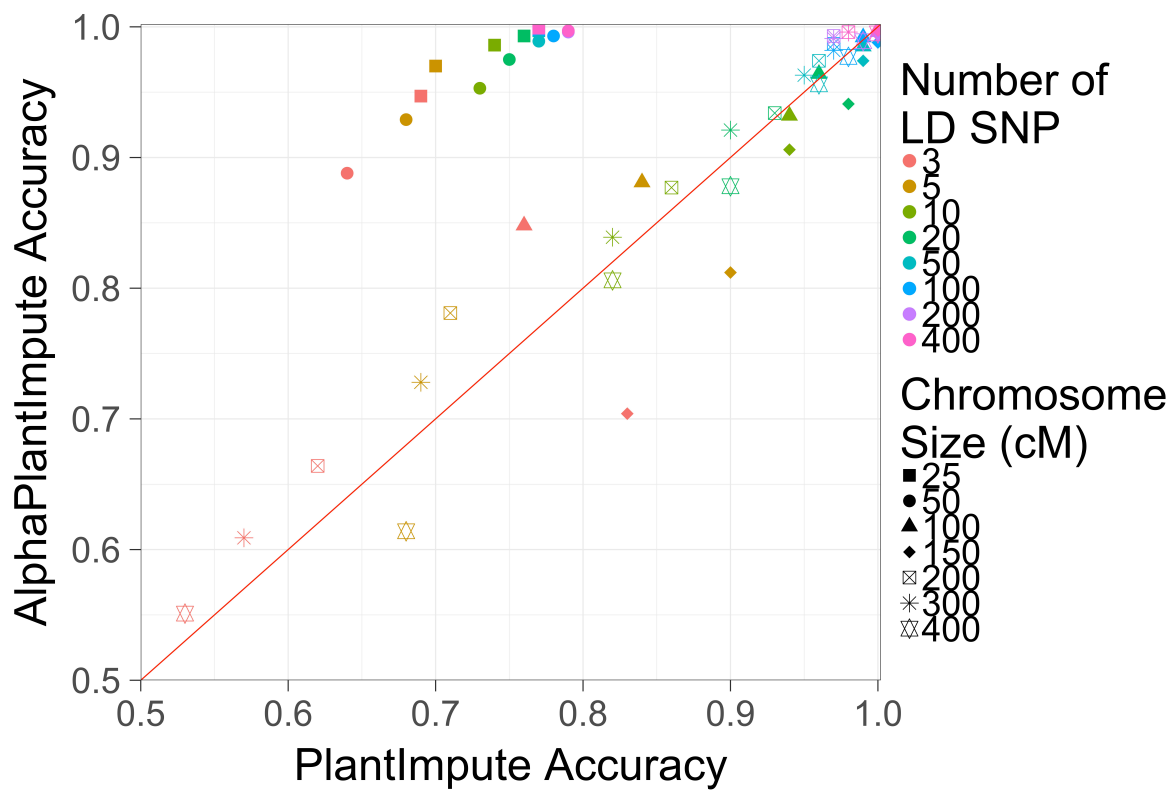
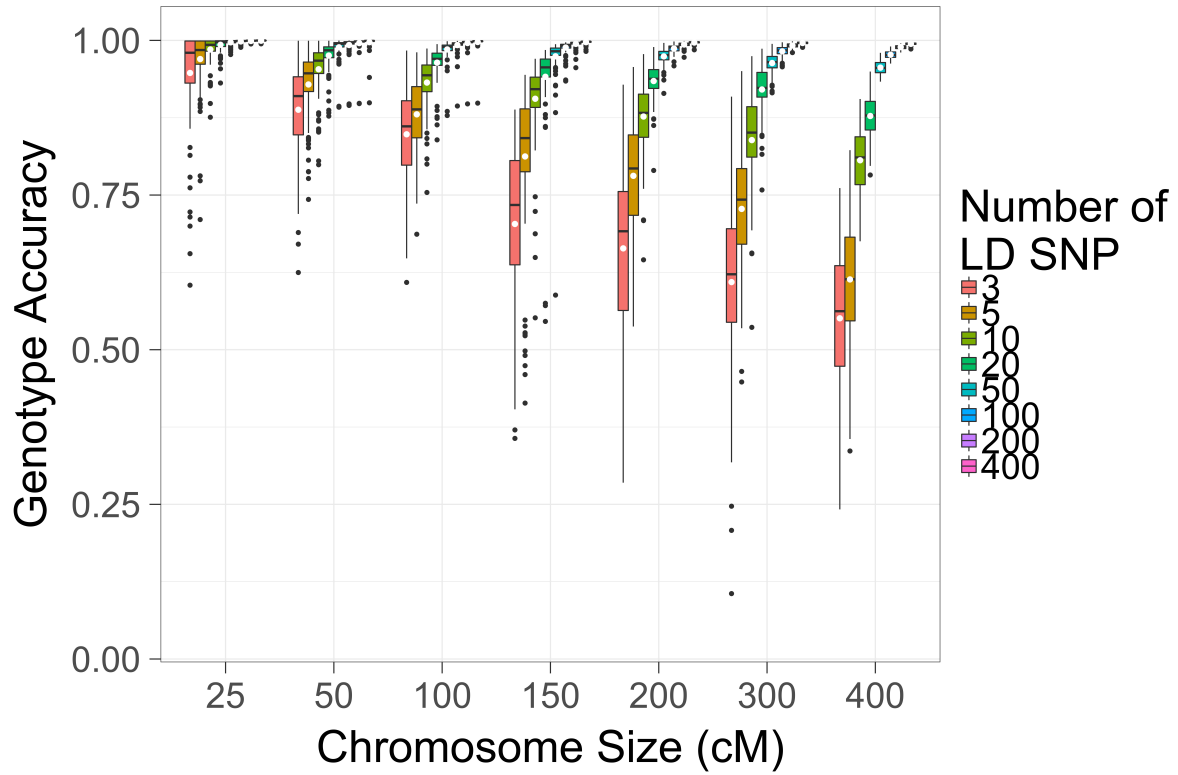


Figure 4 – Effect of the level of inbreeding in parents.

(a) The genotype imputation accuracy using AlphaPlantImpute in F₂ focal individuals of a bi-parental cross where the parents are F₁, F₂, F₄, F₁₀ or F₂₀.

(b) Comparison of the average genotype imputation accuracy using AlphaPlantImpute (y-axis) vs. using PlantImpute (x-axis). The colours represent the different LD panels. The shapes represent the level of inbreeding in the parents. The red diagonal line indicates when the accuracy of PlantImpute equals AlphaPlantImpute. Points above the line are when imputation accuracy is higher with AlphaPlantImpute and points below the line are when imputation accuracy is higher with PlantImpute.

(c) Comparison of the precision in imputation accuracy using AlphaPlantImpute (y-axis) vs. using PlantImpute (x-axis). The colours represent the different LD panels. The shapes represent the level of inbreeding in the parents. The red diagonal line indicates when the precision of PlantImpute equals AlphaPlantImpute. Points above the line indicate when the precision in accuracies is higher in AlphaPlantImpute and points below the line are when the precision in accuracies is higher in PlantImpute.



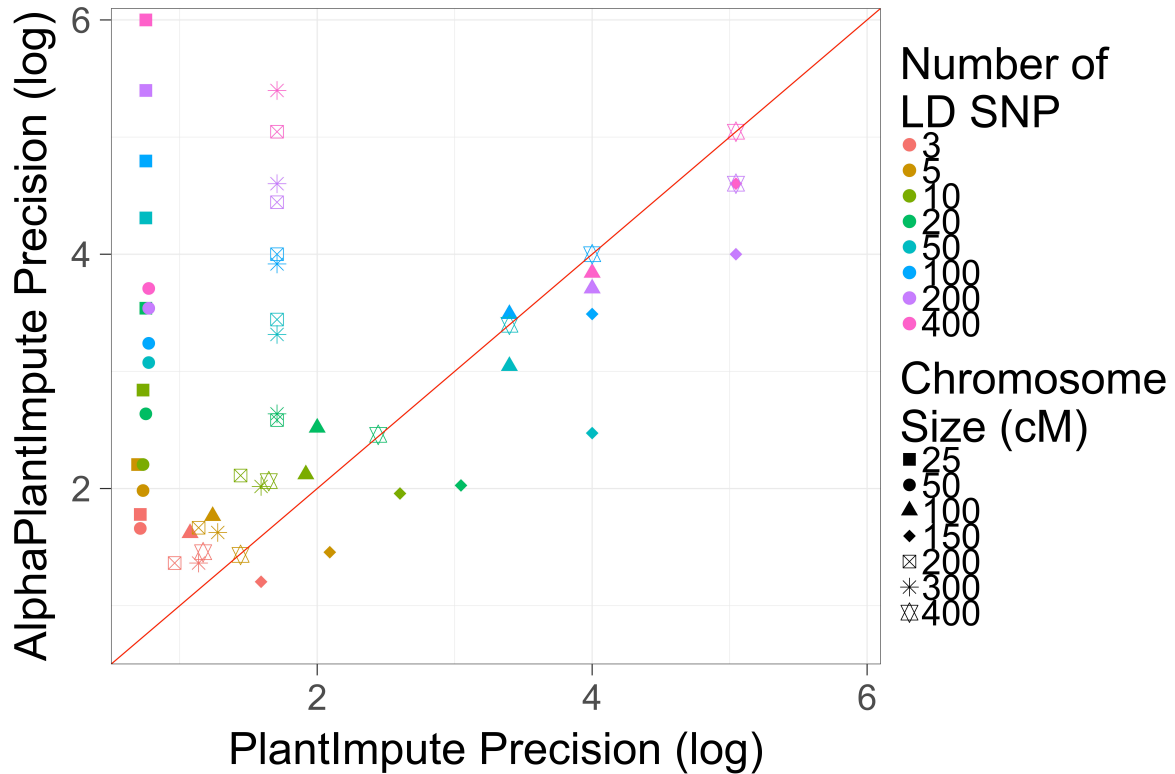


Figure 5 – Effect of chromosome size.

(a) The genotype imputation accuracy using AlphaPlantImpute in F₂ focal individuals from a bi-parental cross of F₂₀ parents against seven chromosome sizes of 25, 50, 100, 150, 200, 300, and 400 cM.

(b) Comparison of the average genotype imputation accuracy using AlphaPlantImpute (y-axis) vs. using PlantImpute (x-axis). The colours represent the different LD panels. The shapes represent the chromosome size. The red diagonal line indicates when the accuracy of PlantImpute equals AlphaPlantImpute. Points above the line are when imputation accuracy is higher with AlphaPlantImpute and points below the line are when imputation accuracy is higher with PlantImpute.

(c) Comparison of the precision in imputation accuracy using AlphaPlantImpute (y-axis) vs. using PlantImpute (x-axis). The colours represent the different LD panels. The shapes represent the chromosome size. The red diagonal line indicates when precision of PlantImpute equals AlphaPlantImpute. Points above the line indicate when the precision in accuracies is higher in AlphaPlantImpute and points below the line are when the precision in accuracies is higher in PlantImpute.

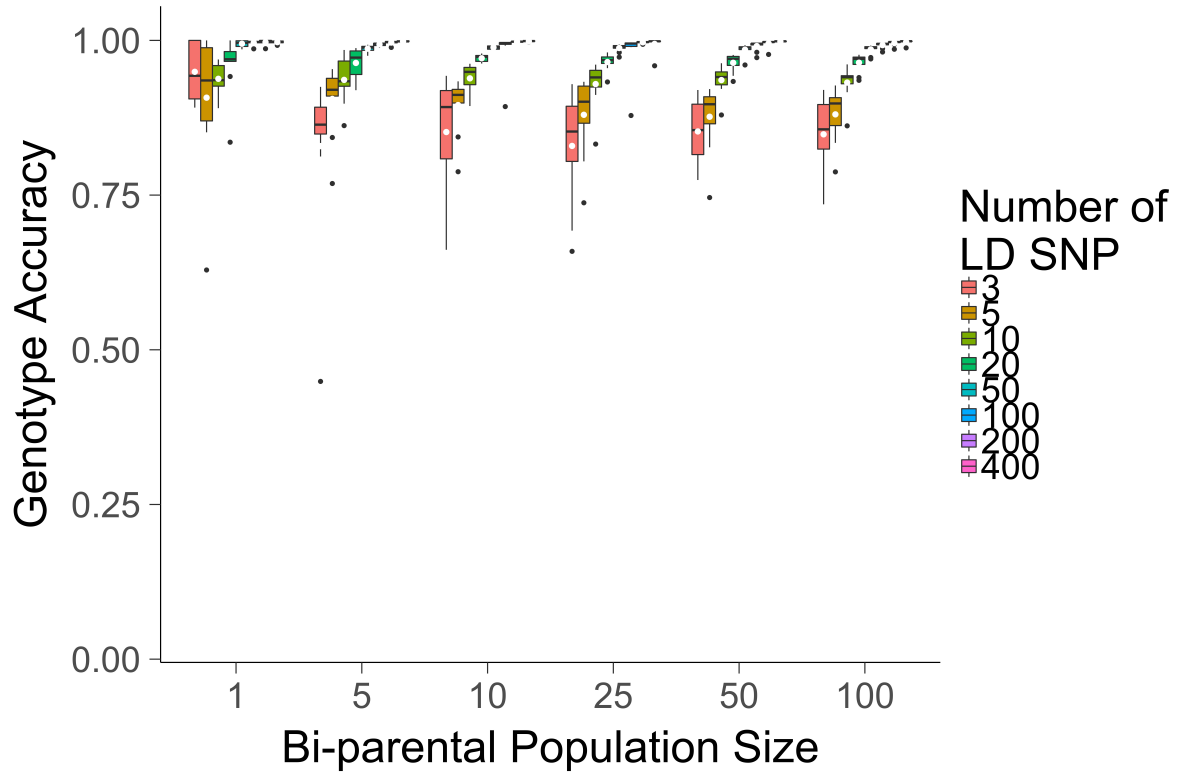


Figure 6 – Effect of the number of focal individuals in the bi-parental population.

The number of focal individuals in the bi-parental population against the genotype imputation accuracy using AlphaPlantImpute for F_2 focal individuals of a bi-parental cross where the parents are F_{20} inbred individuals.

Table 1 – Computational requirements of AlphaPlantImpute

Parents	Focal Individuals	LD panel	Time (Seconds)	Memory (Gb)
F20	F2	3	37.41	0.079
F20	F2	50	8.14	0.080
F20	F2	400	7.95	0.082
F20	F10	3	49.33	0.079
F20	F10	50	8.48	0.080
F20	F10	400	9.40	0.082
F1	F2	3	26.70	0.080
F1	F2	50	35.10	0.080
F1	F2	400	12.58	0.082
F1	F10	3	24.66	0.079
F1	F10	50	35.41	0.080
F1	F10	400	10.35	0.082
Average:			22.13	0.080

