1

2

3

4       Assessment of the performance of different hidden Markov models for imputation in animal

5                                              breeding

6           Andrew Whalen, Gregor Gorjanc, Roger Ros-Freixedes, and John M Hickey

7

8         The Roslin Institute and Royal (Dick) School of Veterinary Studies, The University of

9                               Edinburgh, Midlothian, Scotland, UK

10

11                                              Abstract

12    In this paper we review the performance of various hidden Markov model-based imputation

13    methods in animal breeding populations. Traditionally, heuristic-based imputation methods have

14    been used for imputation in large animal populations due to their computational efficiency,

15    scalability, and accuracy. However, recent advances in the area of human genetics have

16    increased the ability of probabilistic hidden Markov model methods to perform accurate phasing

17    and imputation in large populations. These advances may enable these methods to be useful for

18    routine use in large animal populations. To test this, we evaluate here the accuracy and

19    computational cost of several methods in a series of simulated populations and a real animal

20    population. We first tested single-step (diploid) imputation, which performs both phasing and

21    imputation. Then we tested pre-phasing followed by haploid imputation. We tested four diploid

22    imputation methods (fastPHASE, Beagle v4.0, IMPUTE2, and MaCH), three phasing methods,

23    (SHAPEIT2, HAPI-UR, and Eagle2), and three haploid imputation methods (IMPUTE2, Beagle

24    v4.1, and minimac3). We found that performing pre-phasing and haploid imputation was faster

25    and more accurate than diploid imputation. In particular, we found that pre-phasing with Eagle2

26    or HAPI-UR and imputing with minimac3 or IMPUTE2 gave the highest accuracies in both

27    simulated and real data.

28 **Introduction**

29  In this paper we review and analyse the use of hidden Markov model (HMM) based

30 imputation methods for animal breeding populations. Genotype imputation is a key aspect of

31 many modern animal breeding programs and allows genetic information to be obtained on a

32 large number of animals at a low cost. When imputation is applied to a breeding program, a

33 small subset of individuals (e.g., sires) are genotyped at high density, and the remaining animals

34 are genotyped at a lower density. Statistical regularities between shared chromosomal segments

35 are used to fill in the untyped loci. Modern imputation methods fill in missing genotypes at a

36 very high accuracy (e.g., Hickey et al., 2012; Sargolzaei et al., 2011), increasing the number of

37 animals that can be genotyped for a fixed budget. The larger pool of genotyped animals increases

38 the accuracy of genetic predictions on all animals (Daetwyler et al., 2008) and offers the

39 potential to increase selection intensity.

40  Traditionally, heuristic imputation methods have dominated animal breeding (Hickey et al.,

41 2012; Sargolzaei et al., 2011; VanRaden et al., 2013). These heuristic methods use large

42 chromosome segments shared between closely related animals to rapidly and accurately impute

43 untyped or otherwise missing loci. In contrast, imputation methods used in human genetics have

44 largely been based on the probabilistic HMM framework of Li and Stephens (2003). These

45 probabilistic methods tend to have higher accuracy than heuristic methods in datasets where

46 individuals are not closely related. However, these methods have come at too high of a

47 computational cost for routine imputation in animal populations.

48  In the last few years, the speed of HMM methods has improved. They have been used to

49 impute hundreds of thousands of individuals to hundreds of thousands of loci in reasonable

50 computational time (Browning and Browning, 2016; Loh et al., 2016a). These improvements

51    have been driven by the widespread availability of large haplotype reference panels, and the

52    emergence of a two-step imputation pipeline where observed genotypes are first phased and then

53    untyped loci are imputed based on their phased haplotypes (Spiliopoulou et al., 2017). The

54    improved scaling of HMMs may allow for their routine use in large animal breeding populations.

55    However, given the lack of appropriate public domain haplotype reference panels for many

56    animal populations, smaller population sizes, and sparser marker density, it is not clear that the

57    advances in HMMs will be realized for animal imputation. Furthermore, there are a number of

58    competing HMM imputation methods and it is not clear which is most suited for routine use in

59    animal breeding.

60        In this paper we provide a high-level review of several imputation methods and study their

61    performance on simulated and real data. We grouped comparisons based on single-step (diploid)

62    imputation methods and a two-step combination of pre-phasing and haploid imputation methods.

63    Specifically, for diploid imputation we test fastPHASE (Scheet and Stephens, 2006), Beagle v4.0

64    (Browning and Browning, 2007), IMPUTE2 (Howie et al., 2009), and MaCH (Li et al., 2010).

65    For pre-phasing we test SHAPEIT2 (Delaneau et al., 2012), HAPI-UR (Williams et al., 2012),

66    and Eagle2 (Loh et al., 2016b), followed by haploid imputation with IMPUTE2 (Howie et al.,

67    2009), Beagle v4.1 (Browning and Browning, 2016), or minimac3 (Das et al., 2016). We first

68    review these methods and then evaluate the performance of these methods on simulated and real

69    data.

70                                    **Hidden Markov Models**

71        All of the methods considered are based on Li and Stephens' (2003) HMM framework.

72    Under this framework an individual's genotype is considered to be a mosaic of haplotypes from

73    a haplotype reference panel $H=\{h_1...h_K\}$. The methods calculate the probability that the

74   individual has the pair of haplotypes, $h_j$ and $h_k$ at a locus $i$ given the observed genotype ($g_i$),

75   $p(h_{ij}, h_{ik}/g_i)$. To account for linkage between adjacent loci, the methods evaluate the probability of

76   a haplotype based on its fit to the observed genotypes at the loci and its similarity to the

77   haplotypes inferred at nearby loci:

78   $$p(h_{ij}, h_{ik}/\, g) = p(h_{ij}, h_{ik}/g_i)p(h_{ij}, h_{ik}/h_{i-1}, h_{i+1})p(h_{i-1}/g_{-i})p(h_{i+1}/g_{+i}). \qquad (1)$$

79   The term $p(h_{ij}, h_{ik}/g_i)$ measures the fit between the pair of haplotypes and the observed genotype

80   at a locus. The term $p(h_{ij}, h_{ik}/h_{i-1}, h_{i+1})$ captures transitions between haplotypes given the

81   haplotypes at neighbouring loci. The terms $p(h_{i-1}/g_{-i})$ and $p(h_{i+1}/g_{+i})$ measure the fit between

82   haplotypes and observed genotypes at the remaining loci. These probabilities can be calculated

83   using the standard forward-backward algorithm (Rabiner, 1989).

84           Traditionally, methods that rely on the Li and Stephens framework scale linearly with

85   both the number of individuals and the number of loci and quadratically with the number of

86   reference haplotypes. The quadratic scaling is due to phase uncertainty at heterozygous loci,

87   requiring the methods to model haplotypes assigned on both chromosomes simultaneously. The

88   quadratic scaling quickly leads to intractable computational costs even for small reference

89   panels, but can be avoided if the low-density individuals are pre-phased, which allows

90   haplotypes to be considered independently. Haploid imputation, imputation with pre-phased

91   haplotypes, therefore scales linearly with the number of individuals, number of loci, and number

92   of reference haplotypes.

93           In this paper we consider two classes of HMMs. In the first class, diploid imputation

94   methods perform phasing and imputation simultaneously, resulting in quadratic scaling with the

95   reference panel size. To mitigate this issue, each of the evaluated methods, fastPHASE, Beagle

96   v4.0, IMPUTE2, and MACH, employ their own strategy to reduce the effective number of

97    reference haplotypes while maintaining high accuracy. In contrast, two-step imputation methods

98    treat phasing and imputation as separate problems. Individuals are first phased and then imputed

99    using a haploid HMM which scales linearly with the number of reference haplotypes. Phasing

100    methods may have either quadratic, super-linear, or linear dependence on the number of

101    reference haplotypes. A number of tricks are deployed to increase phasing speed and accuracy

102    that would not be applicable if the phasing methods also needed to handle genotype uncertainty

103    at untyped loci.

104        Intuitively, we might expect that the diploid imputation methods will have higher

105    accuracy (at a higher computational cost) than separately performing phasing and imputation

106    because they automatically handle phase uncertainty. This is not necessarily the case if most

107    errors in imputation stem from the inability to find appropriate reference haplotypes that would

108    explain observed genotypes.  By performing pre-phasing and then imputation, it may be possible

109    to consider a much larger number of reference haplotypes and thereby increase accuracy by

110    finding a more appropriate set of reference haplotypes which offset accuracy losses due to

111    phasing errors.

112        Below we review methods for diploid imputation, haploid imputation, and phasing.

113    **Diploid imputation**

114        All four diploid imputation methods utilize a haplotype state-space reduction technique to

115    alleviate the impact of modelling a large number of haplotype reference panels. IMPUTE2 and

116    MaCH use subsampling, where the haplotypes considered in each iteration are a sample of the

117    total haplotype pool. fastPHASE and Beagle v4.0 use haplotype clustering, where the overall

118    number of haplotypes is collapsed into a smaller number of "ancestral" haplotypes.

119    In the case of IMPUTE2 and MaCH, each method is run over a series of iterations, and at

120    each iteration a subset of the haplotype reference panel is used to phase and impute individual's

121    genotypes. In MaCH, the subset is selected randomly. In IMPUTE2, the subset is selected to be

122    made up of haplotypes that are "nearby" the currently estimated haplotype for the individual. If

123    these methods are run without an external reference panel, a reference panel is built up from the

124    current phasing of high-density individuals. At each iteration, a new subset of the reference panel

125    is selected for each individual, individuals are imputed and phased based on that subuset, and

126    then a reference panel is re-computed from the currently inferred haplotypes. The methods are

127    run for a small number of iterations (e.g., 20) and the imputation results are averaged across

128    iterations. There is a potential danger in applying these methods in populations of many closely

129    related individuals, due to the potential for feedback between the phasing of closely related

130    relatives (Nettelblad, 2013).

131    In contrast, in fastPHASE and Beagle v4.0 individuals are imputed based on a set of

132    estimated "ancestral" haplotypes. In fastPHASE, an expectation-maximisation (EM) algorithm is

133    used to infer a small number of ancestral haplotypes from the data (e.g., 30) and then iterates

134    between estimating the haplotypes of each individual as a mosaic of ancestral haplotypes, and

135    estimating the ancestral haplotypes based on the haplotype assignments of each individual.

136    Beagle v4.0 uses a similar approach as fastPHASE, but instead of using a fixed number of

137    ancestral haplotypes, it infers the number of ancestral haplotypes at each marker and models the

138    transition between ancestral states at adjacent markers in the form of a directed acyclic graph.

139    **Haploid imputation**

140    In contrast to the four diploid methods, haploid methods do not need to use a state-space

141    reduction technique to handle moderate numbers of haplotypes, because they consider each

142    phased chromosome independently and scale linearly with the number of haplotypes in the

143    reference panels. However, with the recent focus of imputing large bio-bank size human

144    populations (over 100,000 individuals) to whole genome sequence level data, many of the

145    current haploid methods utilize techniques to reduce the computational burden when analyzing

146    large numbers of individuals at a large number of markers.

147          The haploid HMM used by Impute2 is a straightforward extension of the diploid method

148    implemented in the same program. It uses a subset of haplotypes (based on their similarity to the

149    individual's current phasing) to impute individuals. Minimac3 uses a similar technique, but

150    instead of subsetting the reference panel it uses a loss-less haplotype compression technique that

151    combines haplotypes that are identical in a region and updates the likelihood of those haplotypes

152    simultaneously. This update is particularly useful for whole genome sequence data where there

153    may be limited haplotype variation over long windows. Beagle v4.1 moves away from the graph-

154    based haplotype model in Beagle v4.0 and uses a more traditional Li and Stephens model. To

155    reduce computational burden, Beagle v4.1 aggregates adjacent loci together into strings and

156    performs updates based on strings instead of individual markers. In addition it only updates the

157    haplotype probabilities at genotyped loci and linearly interpolates the haplotype probabilities at

158    untyped loci.

159    **Pre-phasing methods**

160          Just as with diploid imputation, HMM-based phasing methods naively scale quadratically

161    with the number of haplotypes in the reference panel. However, this quadratic scaling can be

162    avoided by a state-space reduction technique of splitting the chromosomes into small windows,

163    and assuming that linkage information decays quickly across the window boundaries. Both

164    SHAPEIT2 and HAPI-UR utilize a window-based approach, whereas Eagle2 manages the

165    quadratic dependence by performing a limited beam search through the haplotype space.

166        SHAPEIT2 operates by splitting the chromosome into small haplotype windows, each

167    containing three heterozygous loci. For each window, there are $2^3=8$ possible ways to phase it,

168    and there are $2^6=64$ possible transitions between windows. SHAPEIT2 evaluates the probability

169    of each of the 8 possible haplotypes and 64 transitions based on a haplotype reference panel, and

170    then phases individuals by sampling haplotypes based on their posterior probabilities. The

171    probability of a haplotype in a given window, and transition between windows can be evaluated

172    in a time that scales linearly with the number of reference haplotypes. As in IMPUTE2,

173    SHAPEIT2 subsets the haplotype reference panel by selecting haplotypes that are nearby the

174    current haplotypes of the individual.

175        The window splitting approach may lead to reduced accuracy in animal breeding

176    populations, where individuals are expected to share long chromosome segments. In SHAPEIT2

177    only the between-window transmission probabilities are modeled, and not the probabilities of the

178    underlying reference haplotypes. This means that haplotype assignment information from a given

179    window is only used to update the next window and is ignored for further windows. This

180    approach limits the amount of long range haplotype information (covering more than 3

181    heterozygous loci) that can be exploited. One solution to this is to increase the size of the

182    windows.

183        HAPI-UR takes a similar approach to SHAPEIT2 in reducing the large state-space, but

184    uses a series of growing windows which allow it to exploit longer shared chromosomal

185    segments. In order to process large windows, HAPI-UR takes advantage of a number of

186    computational tricks to drastically reduce computation time. Unlike most methods that assume a

187    small error rate for observed genotypes (to cover genotyping errors, errors in the reference panel,

188    and mutations from the ancestral state), HAPI-UR sets the probability of all reference haplotypes

189    that disagree with the observed haplotype to 0. This allows the evaluation of which haplotypes fit

190    an individual's chromosome to be re-formulated as a bit-wise set-intersection operation. In

191    addition to this, HAPI-UR uses a structured representation of the reference haplotypes that

192    allows for fast lookups of matching haplotypes, and for each individual creates individual

193    specific diploid HMM, which ignores all haplotypes that disagree with homozygote sites. Instead

194    of using a fixed window size, HAPI-UR uses dynamic windows which start small (4 markers)

195    and grows to a user specified maximum (e.g. 64 markers) allowing the method to capture longer

196    chromosome segments.

197         Eagle2 takes a different approach to phasing individuals by not using a window-based

198    haplotype representation. Instead Eagle2 uses a highly efficient reference haplotype storage

199    method based on the positional Burrows-Wheeler Transform (Durbin, 2014) to allow for looking

200    up consistent haplotype pairs in constant time. Instead of employing a full HMM to evaluate all

201    possible haplotypes, Eagle2 employs a beam search to search through only the most promising

202    paths through the space of all possible haplotype pairs. At each heterozygous locus, these paths

203    branch into two possible sub-paths based on the two phasing options. Low probability paths are

204    pruned or merged to keep the overall number of paths small. To decrease the impact that errors

205    in one part of the genome have on subsequent paths, haplotypes are called after 20 markers

206    allowing for the back-propagation of relevant genetic information while decreasing the potential

207    impact of genotyping errors. Absence of approximate window-based haplotype representation

208    makes Eagle2 particularly appealing for animal populations, where a large number of close

209    relatives share long chromosome segments.

210                                **Materials and Methods**

211         We evaluated the performance of the four diploid imputation methods, fastPHASE,

212   Beagle v4.0, IMPUTE2, and MaCH and the three phasing methods, SHAPEIT2, HAPI-UR, and

213   Eagle2 followed by three haploid imputation methods, IMPUTE2, Beagle v4.1, and minimac3 on

214   a series of simulated datasets and a real dataset.

215         The simulated dataset modelled a cattle population. The population consisted of 5

216   generations of 2,000 animals, genotyped on a single chromosome. Each generation was produced

217   by selecting 100 sires from the previous generation based on their true breeding values and

218   randomly mating them with 1,000 dams. The initial set of haplotypes was sampled using a

219   Markovian Coalescent Simulator (Chen et al., 2009) assuming a single 100-cM long

220   chromosome simulated using a per site mutation rate of $2.5 \times 10^{-8}$, and an effective population

221   size (Ne) that changed over time. Based on estimates for the Holstein cattle population (Villa-

222   Angulo et al., 2009), the Ne was set to 100 in the final generation of simulation and to 1256,

223   4350, and 43 500 at 1000, 10 000, and 100 000 generations ago, with linear changes in between.

224   The simulation of breeding values and progeny's haplotypes were performed using AlphaSim

225   (Faux et al., 2016).

226         In the baseline scenario, a single chromosome was genotyped either with a high-density

227   array of 1,000 SNP (allele frequency greater than 0.01) or with a low-density array of 200 SNP,

228   evenly spaced across the high-density array. All of the sires and 100 dams were genotyped at

229   high density. The remaining animals were genotyped at low density. To test the robustness of

230   each method we independently modified the baseline scenario by varying:

231   • the number of SNP in the low-density array from 5 to 400,

232   • the number of individuals in the population from 200 to 10,000, and

233    • the number of genotyped dams from 0 to 500.

234    • We also considered the case when the first two generations were genotyped on a different

235       high-density array from the next two generations, with either 25, 50, or 75% of SNP

236       overlapping between the two high-density arrays.

237          To compare the methods on a real data set, we performed imputation on 56,607

238    individuals from a commercial pig breeding program. These animals were genotyped either with

239    a high-density array of 60,000 SNP or 80,000 SNP or a low-density array of 15,000 SNP. To

240    estimate imputation accuracy, we selected 500 high-density animals (typed at 60,000 SNPs) and

241    masked them to mimic the pattern of missingness found in the SNP of 500 low-density animals.

242    We restricted imputation to chromosome 1.

243          Accuracy was measured with the correlation between animals' imputed genotypes and

244    their true genotypes for each animal separately and averaged over all animals. We did not assess

245    phase accuracy independent of the resulting imputation accuracy.

246          For the simulated datasets, each method was given 8GB of memory and 24 hours to run.

247    Jobs were terminated if they exceeded the runtime or the memory requirements. Unless

248    otherwise specified, we used the default parameters for each simulation. We tested IMPUTE2

249    using either the default 10-cM windows or the entire chromosome and found that imputing the

250    entire chromosome increased accuracy at the cost of additional computational time. We used 5-

251    cM windows with an overlap of 1 cM for Beagle v4.0 and Beagle v4.1. The real dataset was

252    imputed with only the two-step imputation methods given their high accuracy and low runtimes.

253          In all cases, the high-density individuals and low-density individuals were phased

254    separately. For the case of multiple high-density arrays, we used the "merge_ref_panels" option

255    in IMPUTE2 and phased both high-density arrays separately. Because neither minimac3 or

256    Beagle v4.1 accept multiple high-density arrays, we phased the high-density individuals together

257    and let the phasing method fill in the missing genotypes for high-density individuals.

258                                        **Results**

259    **Accuracy**

260            The performance of diploid imputation methods is given in Figure 1. Among the diploid

261    imputation methods, MaCH performs well in most settings. Its accuracy depends slightly on the

262    number of high-density dams, the number of low-density SNPs, and the overlap between high-

263    density arrays. The performance of fastPHASE was similar to that of MaCH, but performed

264    better when there were a small number of high-density animals or small overlap between high-

265    density arrays.  IMPUTE2 had similar accuracy to MaCH, but performed worse when given a

266    small number of high-density dams, or a small number of individuals, and performed better than

267    MaCH when a large number of high-density dams were given. Beagle v4.0 performed similarly

268    to IMPUTE2, but was less affected by the number of high-density dams and number of

269    individuals.

270            The performance of pre-phasing and haploid imputation methods is given in Figure 2.

271    Among these methods, we found that the combination of Eagle2 and IMPUTE2 gave the highest

272    imputation accuracy. Eagle2 led to the highest downstream imputation accuracy regardless of the

273    imputation method, and led to higher accuracies than any of the diploid imputation methods.

274    SHAPEIT2 led to similar but slightly lower performance than Eagle2. HAPI-UR led to the

275    lowest overall performance. Of the tested haploid imputation methods we found only a small

276    difference between IMPUTE2 and Minimac3, but found that Beaglev4.1 had poor imputation

277    accuracy in all tested scenarios. We re-ran Beagle v4.1 with different-sized windows but did not

278    see a noticeable increase in accuracy. There was no interaction between the choice of phasing

279     method and the choice of imputation method for the overall imputation accuracy with the

280     exception of when multiple high-density arrays were used. In this case the combination of HAPI-

281     UR and minimac3 outperformed the combination of Eagle2 and minimac3.

282     **Run time and memory requirements**

283          The elapsed run time of each method in the baseline scenario is given in Table 1. We

284     found that of the diploid imputation methods, MaCH had the lowest run time followed by Beagle

285     v4.0, fastPHASE, and IMPUTE2. Of the phasing methods, HAPI-UR was the fastest by an order

286     of magnitude, followed by Eagle2 and SHAPEIT2. Of the haploid imputation methods,

287     minimac3 was the fastest followed by Beagle v4.1 and IMPUTE2. The combined run-times of

288     the two-step phasing and imputation methods were all substantially lower than that of the single

289     step methods.

290     **Real Data**

291          The performance on the real dataset was similar and is given in Table 4. The imputation

292     accuracy of Eagle2 with minimac3 was 0.992, with Beagle v4.1 was 0.925, and with IMPUTE2

293     was 0.827. The imputation accuracy of HAPI-UR with minimac3 was 0.995%, with Beagle v4.1

294     was 0.939%, and with IMPUTE2 was 0.997%. Phasing with Eagle2 took 7 hours distributed

295     across 8 cores. Phasing with HAPI-UR took 54 hours on a single core. All of the haploid

296     imputation methods took under 6 hours. SHAPEIT2 was not able to phase the high-density and

297     low-density individuals in 4 days and so was not analysed.

298                                        **Discussion**

299          In this paper we evaluated the performance of HMM based imputation methods for

300     imputation in animal populations. We found that combinations of phasing and haploid

301     imputation methods provide increased imputation accuracy at substantially reduced runtimes

302    compared to diploid imputation methods. The combination of using Eagle2 to pre-phase

303    individuals and using minimac3 to impute the data lead to high accuracy imputation in a wide

304    range of simulation scenarios and when analysing a real animal population.

305        The results of this paper highlight the power of separately phasing and imputing

306    individuals. Intuitively it makes sense that performing phasing and imputation in a single step

307    may increase imputation accuracy by marginalizing over uncertainty in phasing. However, the

308    results here suggest that the additional accuracy lost by marginalizing over phasing errors is

309    outweighed by the accuracy gained by considering larger haplotype reference panels. These

310    results are particularly surprising in the context of animal populations where pre-existing

311    reference panels may not exist (at least in the public domain), and so the reference panel itself is

312    inferred by phasing high-density genotyped individuals. Our results suggest that modern phasing

313    methods have a sufficiently high accuracy such that this phasing leads to only a small number of

314    errors.

315        The performance of pre-phasing and haploid imputation is also surprising given the lower

316    density of SNP arrays (both high-density and low-density), and the substantially lower number of

317    overall individuals compared to human studies. We found that pre-phasing and haploid

318    imputation was more effective than the best performing diploid imputation method even for a

319    very small number of low-density markers or, low number of high-density dams, and low

320    numbers of individuals.

321        Of the three phasing methods we tested, using Eagle2 led to the most accurate

322    downstream imputation. This is likely due to the fact that Eagle2 is able to exploit longer

323    segments of shared haplotypes between individuals, which are very common in highly related

324    animal populations. Although Eagle2 led to the highest accuracy, we found that HAPI-UR was

325    an order of magnitude faster for most datasets and resulted in a small decrease in accuracy on the

326    simulated scenarios, but no decrease in accuracy on the real dataset. In their original paper, the

327    authors of HAPI-UR suggest that it may be possible to increase the accuracy of HAPI-UR by

328    running it multiple times with different window start positions and taking the consensus phase

329    (Williams et al., 2012). Due to the low run time, this strategy would be feasible in animal

330    populations but was not analysed here. SHAPEIT2, the oldest of the phasing methods had both

331    the longest run-time which prevented us from evaluating it on the real dataset. Although the

332    authors of SHAPEIT2 have now released SHAPEIT3, they do not recommend using it for

333    populations of under 60,000 individuals and so the performance of SHAPEIT3 was not analysed

334    here.

335        We found little difference in the performance of the assessed haploid imputation

336    methods. Both Minimac3 and IMPUTE2 lead to accurate imputation. The accuracy of IMPUTE2

337    was consistently slightly (<1%) higher than that of minimac3 in simulated data, but the runtime

338    was between two and three times that of minimac3. On the real dataset, the imputation accuracy

339    of IMPUTE2 dropped when Eagle2 was used to pre-phase the data, but remained high when

340    HAPI-UR was used to pre-phase the data. Overall the performance of Beagle v4.1 was poor for

341    performing haploid imputation, although improved when analysing the real data set. This may be

342    a result of the approximations used in Beagle v4.1, which were designed for imputation of

343    human high-density SNP arrays to whole genome sequence data. These approximations seem

344    less appropriate for low-density SNP arrays used in some animal populations.

345        With two exceptions, we found little interaction between the choice of phasing method

346    and the choice of haploid imputation method. The first exception came in the performance of

347    HAPI-UR when individuals were genotyped with multiple, semi-overlapping, SNP arrays. In this

348   case the performance of HAPI-UR with minimac3 or Beagle v4.1 was substantially higher than

349   the performance of Eagle2 with minimac3 or Beagle v4.1, although the accuracy of HAPI-UR

350   with IMPUTE2 remained lower than that of Eagle2 with IMPUTE2. The underlying reason for

351   this difference stems from the fact that in the case of minimac3 and Beagle v4.1 the phasing

352   algorithms were also used to perform imputation on the missing non-overlapping SNPs in each

353   high-density array, whereas in IMPUTE2 the two high-density arrays were phased separately,

354   and IMPUTE2 was used to fill in missing SNPs as part of it's high-density array merging step.

355   The increased accuracy with HAPI-UR over Eagle2 in this scenario suggests that HAPI-UR can

356   impute untyped loci in high-density arrays better than Eagle2. This is consistent with the second

357   exception where HAPI-UR led to as high imputation accuracy, if not higher, as Eagle2 when

358   performing imputation on the real dataset. Animals in the real dataset were genotyped with two

359   high-density arrays, and two low-density arrays, and also exhibited a number of spontaneously

360   missing SNPs. When using Eagle2 to phase individuals, IMPUTE2 and Beagle v4.1 markedly

361   decreased in performance, particularly compared to minimac3. In contrast when HAPI-UR was

362   used to phase individuals the performance of minimac3, IMPUTE2 and Beagle v4.1 remained

363   high, suggesting an advantage of using HAPI-UR over Eagle2 when individuals are genotyped

364   on multiple arrays or when observing a large amount of spontaneous missingness.

365       Some of the analysed phasing methods have an option to use pedigree information to

366   improve phasing. Although these options were originally designed to help phase and impute

367   parent-progeny trios (Browning and Browning, 2009), they can also be used for larger pedigrees

368   (O'Connell et al., 2014). Previous work in phasing and imputing animal populations has found

369   that combining pedigree and linkage information can improve phasing and imputation accuracy

370   (Hickey et al., 2012). In this paper, we did not analyse the option to use pedigree information,

371    but focused solely on HMMs based methods that use linkage-disequilibrium information for

372    phasing and imputation as originally proposed by Li and Stephens (2003). SHAPEIT2

373    (O'Connell et al., 2014), Beagle v4.0 (Browning and Browning, 2009), and HAPI-UR (Williams

374    et al., 2012) all provide options to use parent-progeny trio information. However, the two top

375    performing methods, Eagle2 and minimac3, do not provide this option. Future work is needed to

376    analyse how HMMs can utilize pedigree information to improve phasing and imputation, and to

377    merge these insights with high-performance methods reviewed and tested here.

378        Overall, this study suggests that modern pre-phasing and haploid imputation methods can

379    perform fast and accurate imputation of animal populations of any size. We noticed no

380    disadvantage of using the two-step imputation approach even in cases of small populations, low-

381    density SNP arrays, or multiple high-density arrays. Of the algorithms, we found that Eagle2 and

382    HAPI-UR both reliably pre-phased the data and that IMPUTE2 and minimac3 lead to the highest

383    imputation accuracy. However, we also noted a decreased accuracy when Eagle2 and IMPUTE2

384    were used to pre-phase and impute the data when animals were genotyped with semi-overlapping

385    high-density SNP arrays. In this case the usage of Eagle 2 with minimac3 and HAPI-UR with

386    IMPUTE2 or minimac3 lead to high accuracy. Overall, the results of these studies highlight the

387    importance and feasibility of using HMMs to perform imputation in animal populations even as

388    an increasing number of animals are genotyped and as genotyping densities increase.

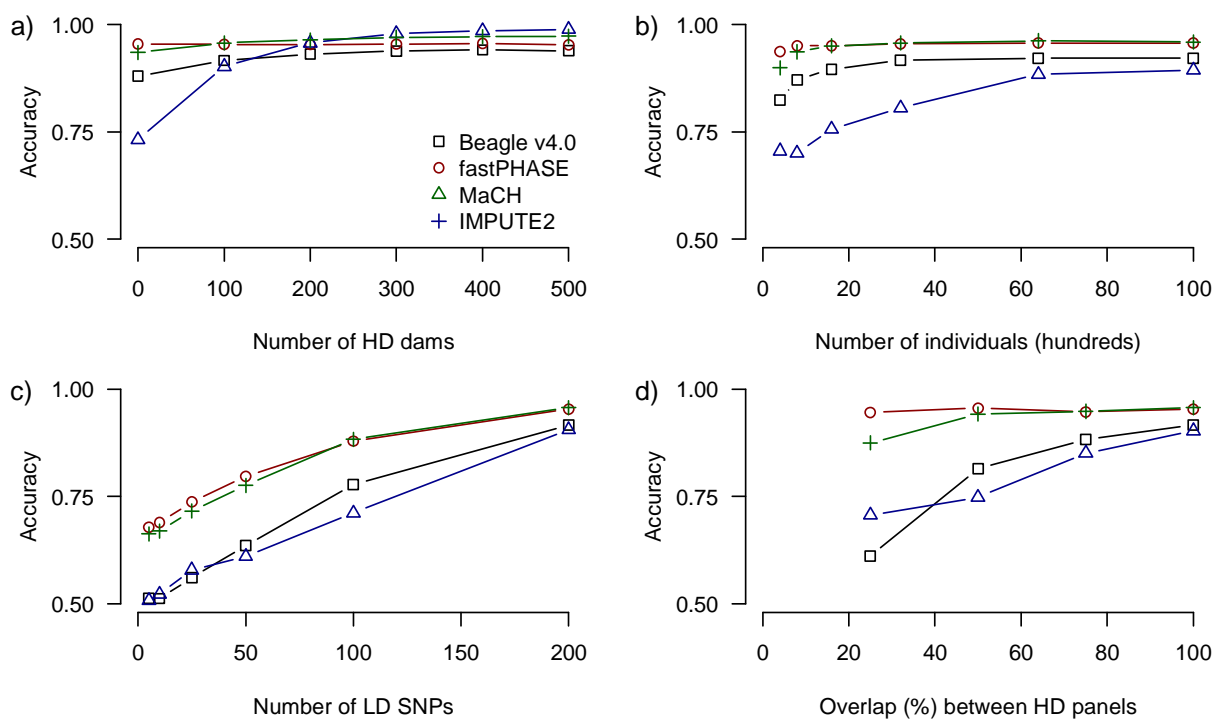389                                **Acknowledgements**

395 **References**

396    Browning, B.L., and Browning, S.R. (2009). A Unified Approach to Genotype Imputation and
397    Haplotype-Phase Inference for Large Data Sets of Trios and Unrelated Individuals. Am. J. Hum.
398    Genet. *84*, 210–223.

399    Browning, B.L., and Browning, S.R. (2016). Genotype Imputation with Millions of Reference
400    Samples. Am. J. Hum. Genet. *98*, 116–126.

401    Browning, S.R., and Browning, B.L. (2007). Rapid and accurate haplotype phasing and missing-
402    data inference for whole-genome association studies by use of localized haplotype clustering.
403    Am. J. Hum. Genet. *81*, 1084–1097.

404    Chen, G.K., Marjoram, P., and Wall, J.D. (2009). Fast and flexible simulation of DNA sequence
405    data. Genome Res. *19*, 136–142.

406    Daetwyler, H.D., Villanueva, B., and Woolliams, J.A. (2008). Accuracy of Predicting the
407    Genetic Risk of Disease Using a Genome-Wide Approach. PLoS ONE *3*, e3395.

408    Das, S., Forer, L., Schönherr, S., Sidore, C., Locke, A.E., Kwong, A., Vrieze, S.I., Chew, E.Y.,
409    Levy, S., McGue, M., et al. (2016). Next-generation genotype imputation service and methods.
410    Nat. Genet. *advance online publication*.

411    Delaneau, O., Marchini, J., and Zagury, J.-F. (2012). A linear complexity phasing method for
412    thousands of genomes. Nat Meth *9*, 179–181.

413    Durbin, R. (2014). Efficient haplotype matching and storage using the positional Burrows–
414    Wheeler transform (PBWT). Bioinformatics *30*, 1266–1272.

415    Faux, A.-M., Gorjanc, G., Gaynor, R.C., Battagin, M., Edwards, S.M., Wilson, D.L., Hearne,
416    S.J., Gonen, S., and Hickey, J.M. (2016). AlphaSim: Software for Breeding Program Simulation.
417    Plant Genome *9*.

418    Hickey, J.M., Kinghorn, B.P., Tier, B., van der Werf, J.H., and Cleveland, M.A. (2012). A
419    phasing and imputation method for pedigreed populations that results in a single-stage genomic
420    evaluation. Genet. Sel. Evol. *44*, 11.

421    Howie, B.N., Donnelly, P., and Marchini, J. (2009). A flexible and accurate genotype imputation
422    method for the next generation of genome-wide association studies. PLoS Genet. *5*, e1000529.

423    Li, N., and Stephens, M. (2003). Modeling linkage disequilibrium and identifying recombination
424    hotspots using single-nucleotide polymorphism data. Genetics *165*, 2213–2233.

425    Li, Y., Willer, C.J., Ding, J., Scheet, P., and Abecasis, G.R. (2010). MaCH: using sequence and
426    genotype data to estimate haplotypes and unobserved genotypes. Genet. Epidemiol. *34*, 816–834.

427    Loh, P.-R., Danecek, P., Palamara, P.F., Fuchsberger, C., Reshef, Y.A., Finucane, H.K.,
428    Schoenherr, S., Forer, L., McCarthy, S., Abecasis, G.R., et al. (2016a). Reference-based phasing
429    using the Haplotype Reference Consortium panel. Nat. Genet. *48*, 1443–1448.

430    Loh, P.-R., Danecek, P., Palamara, P.F., Fuchsberger, C., A Reshef, Y., K Finucane, H.,
431    Schoenherr, S., Forer, L., McCarthy, S., Abecasis, G.R., et al. (2016b). Reference-based phasing
432    using the Haplotype Reference Consortium panel. Nat. Genet. *48*, 1443–1448.

433    Nettelblad, C. (2013). Breakdown of methods for phasing and imputation in the presence of
434    double genotype sharing. PloS One *8*, e60354.

435    O'Connell, J., Gurdasani, D., Delaneau, O., Pirastu, N., Ulivi, S., Cocca, M., Traglia, M., Huang,
436    J., Huffman, J.E., and Rudan, I. (2014). A general approach for haplotype phasing across the full
437    spectrum of relatedness. PLoS Genet. *10*, e1004234.

438    Rabiner, L. (1989). A tutorial on hidden Markov models and selected applications in speech
439    recognition. Proc. IEEE *77*, 257–286.

440    Sargolzaei, M., Chesnais, J.P., and Schenkel, F.S. (2011). FImpute - An efficient imputation
441    algorithm for dairy cattle populations. J. Dairy Sci. *94 (E-Suppl. 1)*, 421.

442    Scheet, P., and Stephens, M. (2006). A Fast and Flexible Statistical Model for Large-Scale
443    Population Genotype Data: Applications to Inferring Missing Genotypes and Haplotypic Phase.
444    Am. J. Hum. Genet. *78*, 629–644.

445    Spiliopoulou, A., Colombo, M., Orchard, P., Agakov, F., and McKeigue, P. (2017). GeneImp:
446    Fast Imputation to Large Reference Panels Using Genotype Likelihoods from Ultra-Low
447    Coverage Sequencing. Genetics.

448    VanRaden, P.M., Null, D.J., Sargolzaei, M., Wiggans, G.R., Tooker, M.E., Cole, J.B.,
449    Sonstegard, T.S., Connor, E.E., Winters, M., van Kaam, J.B.C.H.M., et al. (2013). Genomic
450    imputation and evaluation using high-density Holstein genotypes. J. Dairy Sci. *96*, 668–678.

451    Villa-Angulo, R., Matukumalli, L.K., Gill, C.A., Choi, J., Tassell, C.P.V., and Grefenstette, J.J.
452    (2009). High-resolution haplotype block structure in the cattle genome. BMC Genet. *10*, 19.

453    Williams, A.L., Patterson, N., Glessner, J., Hakonarson, H., and Reich, D. (2012). Phasing of
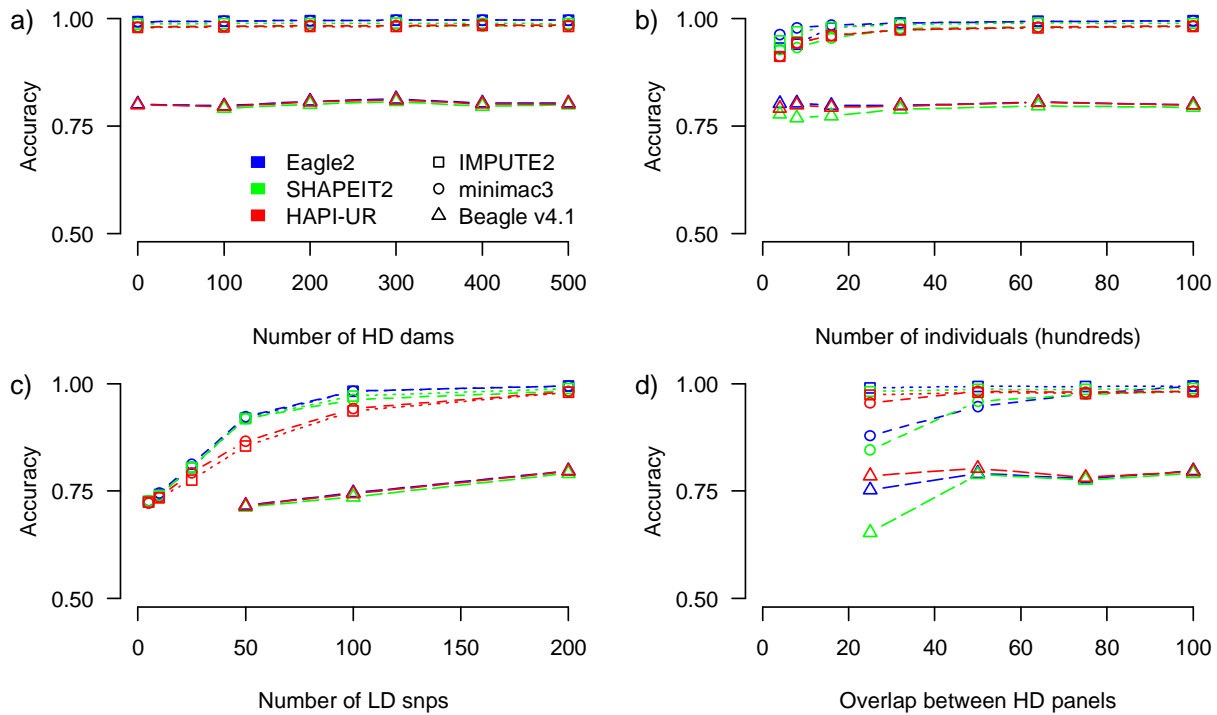454    Many Thousands of Genotyped Samples. Am. J. Hum. Genet. *91*, 238–251.

455

456



457

*Figure 1. Performance of each diploid HMM algorithm for each set of simulations. Unless otherwise noted there were 1000 high-density SNPs, 200 low-density SNPs, 100 dams genotyped at high-density and complete overlap between the high-density arrays of generations 1 and 2 and those of 3 and 4. We varied (a) the number of dams genotyped at high-density, (b) the number of individuals in the population, (c) the number of SNPs in the low-density array, and (d) the amount of overlap between the high-density array for generations 1 and 2 and those of 3 and 4.*

465

466

467

*Figure 2. Performance of each combination of pre-phasing and haploid HMM method. Unless*

*otherwise noted there were 1000 high-density SNPs, 200 low-density SNPs, 100 dams*

*genotyped at high-density and complete overlap between the high-density arrays of generations*

*1 and 2 and those of 3 and 4. We varied (a) the number of dams genotyped at high-density, (b)*

*the number of individuals in the population, (c) the number of SNPs in the low-density array,*

*and (d) the amount of overlap between the high-density array for generations 1 and 2 and those*

*of 3 and 4.*

477      Table 3

478      Simulated data: Run time and accuracy for diploid imputation, phasing, and haploid imputation
479      methods in the baseline scenario. The run time is given in seconds separately for phasing and
480      imputation steps and as a total.
481

| Phasing method | Imputation method | HD Phasing (s) | LD Phasing (s) | Imputation (s) | Total (s) | Accuracy |
|---|---|---|---|---|---|---|
| / | IMPUTE2 | / | / | 42,796 | 42,796 | 0.861 |
| / | Beagle v4.0 | / | / | 23,042 | 23,042 | 0.901 |
| / | MaCH | / | / | 21,998 | 21,998 | 0.944 |
| / | fastPHASE | / | / | 28,892 | 28,892 | 0.941 |
| HAPI-UR | IMPUTE2 | 117 | 14 | 149 | 280 | 0.964 |
| HAPI-UR | minimac3 | 117 | 14 | 62 | 193 | 0.967 |
| HAPI-UR | Beagle v4.1 | 117 | 14 | 78 | 209 | 0.793 |
| Eagle2 | IMPUTE2 | 1,361 | 207 | 148 | 1,717 | 0.988 |
| Eagle2 | minimac3 | 1,361 | 207 | 55 | 1,623 | 0.988 |
| Eagle2 | Beagle v4.1 | 1,361 | 207 | 79 | 1,647 | 0.794 |
| SHAPEIT2 | IMPUTE2 | 8,495 | 1,175 | 150 | 9,820 | 0.979 |
| SHAPEIT2 | minimac3 | 8,495 | 1,175 | 58 | 9,728 | 0.977 |
| SHAPEIT2 | Beagle v4.1 | 8,495 | 1,175 | 77 | 9,747 | 0.792 |

482

483

484    Table 4

485    Real data: Run time and accuracy for phasing, and haploid imputation methods on the real
486    dataset scenario. The run time is given in hours separately for phasing and imputation steps and
487    as a total. For Eagle2, the program was run distributed across 8 compute cores. HAPI-UR was
488    run on a single core.
489

| Phasing method | Imputation method | HD Phasing (h) | LD Phasing (h) | Imputation (h) | Total (h) | Accuracy |
|---|---|---|---|---|---|---|
| HAPI-UR | IMPUTE2 | 11.53 | 43.09 | 60.25 | 12.48 | 0.997 |
| HAPI-UR | minimac3 | 11.53 | 43.09 | 56.89 | 9.06 | 0.995 |
| HAPI-UR | Beagle v4.1 | 11.53 | 43.09 | 57.32 | 11.04 | 0.939 |
| Eagle2 | IMPUTE2 | 4.48 (8 cores) | 2.37 (8 cores) | 5.63 | 12.48 | 0.827 |
| Eagle2 | minimac3 | 4.48 (8 cores) | 2.37 (8 cores) | 2.21 | 9.06 | 0.992 |
| Eagle2 | Beagle v4.1 | 4.48 (8 cores) | 2.37 (8 cores) | 4.19 | 11.04 | 0.925 |

490